# Scaling up Kernels in 3D CNNs

**Yukang Chen**[1], **Jianhui Liu**[2], **Xiaojuan Qi**[2], **Xiangyu Zhang**[3], **Jian Sun**[3], **Jiaya Jia**[1]

[1]The Chinese University of Hong Kong  [2]The University of Hong Kong  [3]MEGVII Technology

## Abstract

Recent advances in 2D CNNs and vision transformers (ViTs) reveal that large kernels are essential for enough receptive fields and high performance. Inspired by this literature, we examine the feasibility and challenges of 3D large-kernel designs. We demonstrate that applying large convolutional kernels in 3D CNNs has more difficulties in both performance and efficiency. Existing techniques that work well in 2D CNNs are ineffective in 3D networks, including the popular depth-wise convolutions. To overcome these obstacles, we present the *spatial-wise group convolution* and its large-kernel module (SW-LK block). It avoids the optimization and efficiency issues of naive 3D large kernels. Our large-kernel 3D CNN network, *i.e.*, LargeKernel3D, yields non-trivial improvements on various 3D tasks, including semantic segmentation and object detection. Notably, it achieves **73.9%** mIoU on the ScanNetv2 semantic segmentation and **72.8%** NDS nuScenes object detection benchmarks, ranking **1**[st] on the nuScenes LIDAR leaderboard. It is further boosted to **74.2%** NDS with a simple multi-modal fusion. LargeKernel3D attains comparable or superior results than its CNN and transformer counterparts. For the first time, we show that large kernels are feasible and essential for 3D networks. Code and models will be available at github.com/dvlab-research/LargeKernel3D.

## 1 Introduction

3D Sparse convolutional neural networks (CNNs) are widely-used feature extractors in 3D tasks, e.g., semantic segmentation [10, 20] and object detection [52, 13, 63]. The advantages of efficiency and convenient usage ensure its popularity in various applications, such as autonomous driving and robotics. However, 3D sparse CNNs have been recently challenged by transformer-based methods [41, 67, 43], mainly from the view of building up receptive fields. Both global and local [18, 41] self-attention mechanisms are able to capture context information from a large spatial scope. Some parallel literature in 2D Vision Transformers (ViTs) also emphasize their advantages in modeling long-range dependencies [17, 36, 48]. In contrast, common 3D sparse CNNs have limitations in this aspect. The receptive fields of default 3D sparse CNN are constrained by small kernel sizes and spatial disconnection of sparse features (due to the property of submanifold sparse convolution [21]).

Literature in 2D CNNs [38, 15, 55] has already presented a series of methods, combined with large kernels, to enlarge the receptive fields and model capacity. These works present the potential of large-kernel CNNs in 2D vision tasks, but not direct solutions in 3D tasks. ConvNeXt [38] employs 7×7 depth-wise convolution as a strong design, combing with other training techniques to challenge its Swin Transformer counterpart [36]. RepLKNet [15] pursues extremely larger kernel sizes, *e.g.*, 31×31, to boost the performance of downstream tasks. To ensure the effectiveness of RepLKNet [15], additional factors, including depth-wise convolution, are also required. Other work [22] also emphasizes the importance of depth-wise convolution. Due to differences between 3D and 2D tasks, we find that existing 2D design guidelines are infeasible in 3D sparse CNNs.

---

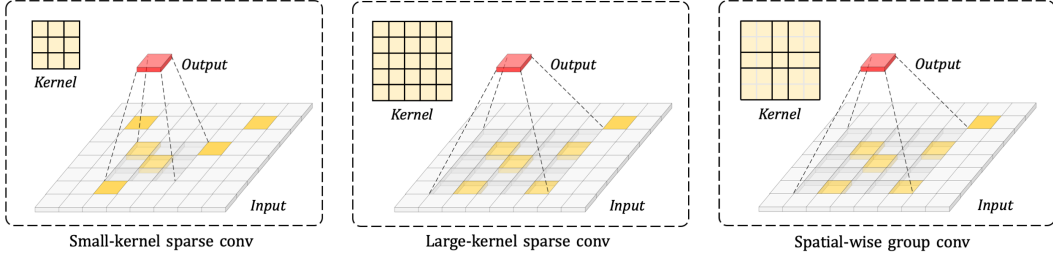Yukang's work was done during an internship in MEGVII Technology.

Figure 1: Process of sparse convolutions with different kernels. *Small-kernel sparse convolution* gathers features in a small local area. It maintains efficiency but discards enough information flow, due to feature disconnection and the small scope. *Large-kernel sparse convolution* is capable of capturing long-range information, at the price of huge amounts of parameters and computations. The proposed *spatial-wise group convolution* remains a spatially large kernel size, but shares weights among local neighbors to keep efficiency. This figure is illustrated in 2D features for simplification.

We summarize the difficulties of 3D large-kernel CNN designs in two aspects. The first obstacle is *efficiency*. 3D convolution has the cubic model size and computations increase ratio, instead of quadratically in 2D convolution. For example, the model size increase by more than $10\times$, when increasing kernels from $3^3$ to $7^3$. The second difficulty exists in *optimization*. In terms of *data amount*, 3D datasets commonly only contain thousands of scenes, which can not match 2D large-scale image benchmarks [12, 34]. Thus, it might not be sufficient to optimize the dramatically proliferated parameters. In terms of *data structure*, different from 2D dense images, 3D sparse data, *e.g.*, point clouds, or voxels, are spatially non-uniform. Not all weights are activated and updated during training. This issue becomes more severe as kernel size is larger, as shown in Fig. 1.

2D literature relies on depth-wise convolutions to train large-kernel networks [15, 38]. It shares weights among channels, making the parameter learning efficient. However, we find this technique, including others working well in 2D CNNs, is ineffective or even harmful in 3D cases (as in Tab. 1). This leads to the question that whether there is a specific solution to 3D large kernels.

To answer this question, we propose spatial-wise group convolution (*SW Conv*) as a solution to the 3D large-kernel design. It is a new family of group convolution by sharing weights among spatially adjacent locations, rather than *depth-wise convolution* [26] of channel-level groups. As shown in Fig. 1, spatial-wise convolution remaps a large kernel (*e.g.*, $7 \times 7$) as a small one (*e.g.*, $3 \times 3$) via grouping spatial neighbors, while the absolutely large spatial size remain unchanged. In terms of the *efficiency* issue, SW Conv occupies less model size for parameter efficiency, identical to small-kernel models. Moreover, it saves actual runtime based on our implementation, compared with plain large kernel counterparts (as in Tab. 3). In terms of the *optimization* problem, weight-sharing among spatial dimensions makes parameters more chance to be updated against sparse input features. We empirically show that SW Conv behaves stably against kernel enlargements (as in Tab. 2).

To make up the detail-capturing ability of large kernels [15], we further introduce a small parallel branch for SW Conv. It forms the spatial-wise large-kernel convolutional block (*SW-LK block*), as shown in Tab. 5. The proposed SW-LK block can readily replace large spatial-wise convolutions into existing 3D convolutional networks. We build up large-kernel backbone networks LargeKernel3D on existing 3D semantic segmentation [10] and object detection [13, 63] networks. Our method enables non-trivial improvements upon previous state-of-the-art methods [13, 63, 10], with a small model complexity overhead. Extensive experiments validate our effectiveness on large-scale benchmarks, including ScanNetv2 [11], KITTI [19] and nuScenes [3]. For 3D semantic segmentation, our method achieves **73.9%** mIoU on ScanNetv2. For object detection, LargeKernel3D achieves **72.8%** NDS on nuScenes, ranking **1**[st] on the nuScenes LIDAR leaderboard. Without bells and whistles, it further improves to **74.2%** NDS in a simple voxel-wise multi-modal fusion manner. It is comparable or superior to previous 3D convolutional and transformer methods.

We visualize the *Effective Receptive Fields* (ERFs) of plain 3D CNNs and our LargeKernel3D in Fig. 2. It shows that deeper small-kernel networks are also constrained by limited ERFs, as sparse features are spatially disconnected while our large-kernel networks resolve this issue. For the first time, we show that large-kernel CNN designs are feasible and effective in 3D tasks.
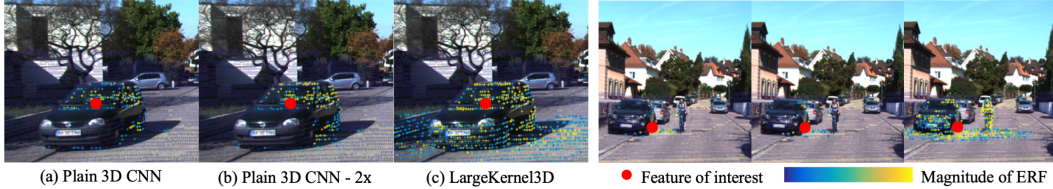
|  (a) Plain 3D CNN | (b) Plain 3D CNN - 2x | (c) LargeKernel3D | ● Feature of interest | Magnitude of ERF |

Figure 2: The *Effective Receptive Fields* (ERF) of plain 3D CNN, plain 3D CNN - 2×, and our LargeKernel3D for 3D object detection. The plain 3D CNN backbone has insufficient ERF size, which hardly covers a nearby car. Plain 3D CNN - 2× is a deeper version with its layers double in each stage. More layers help little in ERF, as the disconnection between sparse features. Our LargeKernel3D obtains large ERF, no matter at the object centers (left) or edges (right), capturing context information. More illustrations are in the *appendix*. It is best viewed in color and by zoom-in.

## 2 Related Work

Large-kernel models have already shown the potential in 2D CNNs, due to the advantages of large receptive fields and high accuracy. We discuss both 2D and 3D related networks in the following.

### 2.1 Large-kernel Models

**Convolutional networks.** Large-kernel settings are not commonly used in 2D convolutional networks until the recent literature [15, 38]. ConvNeXt [38] combines various training techniques and detailed designs, which include large kernel sizes, to compete with Swin Transformer [36]. RepLKNet [15] is more specific work that focuses its main contribution on large-kernel designs. It reveals that the large kernel sizes are more beneficial to downstream tasks, *e.g.*, object detection, and semantic segmentation, rather than image classification. GCNs [44] also presents the improvements from large-kernel designs on semantic segmentation. Other literature approximates large kernel sizes using implicit techniques, including Fourier domain transformation [49], continuous functions [50]. These alternative solutions are not direct and have a large gap from 3D tasks.

**Vision transformer.** Vision transformer networks that conduct attention computation in local areas or windows can also be viewed as large-kernel models. Swin Transformers [36] captures features in shifted windows with sizes from 7 to 12. Its variants [16, 35] present that larger window sizes are beneficial to the performance. Focal Transformer [61] captures fine-grained local attentions with adaptive patch sizes. Inspired by Swin Transfomers [36], SST [18] and Stratified Transformer [29] apply similar window-based self-attention on 3D object detection and semantic segmentation, which hint that the large-kernel models might also be effective in 3D networks.

### 2.2 3D Neural Networks

**LIDAR feature extractors.** A key challenge in 3D tasks is to learn effective representations from the sparse and non-uniformed 3D geometric data, *e.g.*, point clouds. In general, there are two kinds of 3D feature extractors. The first is to learn on point cloud directly, using a series of PointNet networks [46, 47]. These methods are direct solutions, but infeasible for large-scale scenes, where millions of points exist. In addition, the sampling and grouping operations in PointNet++ [47] are also time-consuming. Some follow-ups [67, 43] use grid-sampling or scene-partition to save computation in each forward and combine results from multiple runs. The second is to process point clouds with voxelizations and apply 3D sparse CNNs. Due to its efficiency advantages, this stream of methods has been widely used in various 3D tasks, such as MinkowskiNet [10] in 3D semantic segmentation and sparse 3D CNNs backbone networks in 3D object detectors [52, 13, 63].

**Multi-modal fusion.** Multi-modal fusion is another important technique for 3D feature extraction, especially for 3D object detection. Existing approaches include point/voxel-level [56, 32, 64], object-level [2, 7], and BEV [37, 33] fusion methods. Point/voxel-level fusion methods directly paint image features onto LIDAR points or voxels, instead of relying on Region of Interests (RoIs) prediction or Bird's-Eye View (BEV) encoders. We fuse image features onto voxels at the beginning layer of our 3D backbone, to demonstrate the representation capacity of our large-kernel backbone.
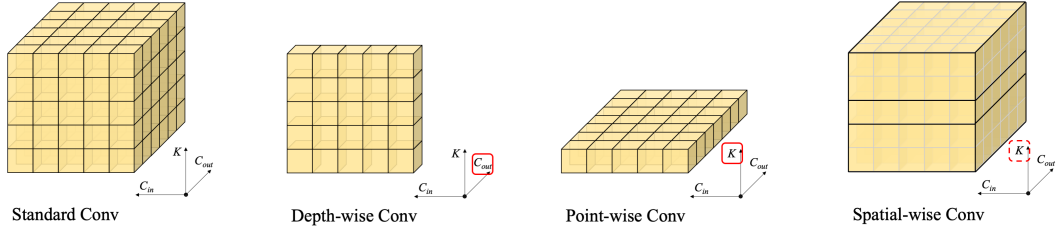
Figure 3: Illustrations on different convolution types. We formulate convolution weight into 3 dimensions, *i.e.*, input channel $C_{in}$, output channel $C_{out}$, and the kernel spatial dimension $K$. Depth-wise convolution splits groups along channels. Point-wise convolution fixes kernel dimension to be 1. Spatial-wise group convolution shares weights among the kernel spatial dimension.

## 3  3D Large-kernel Convolutional Network

### 3.1  Revisiting 3D Sparse CNNs

**Preliminary.** 3D sparse CNNs consist of 3D sparse convolutions, which are typically regular and submanifold sparse convolutions [21]. We formulate 3D sparse convolutions in Eq. 1. Given a set of sparse input features $\{\mathrm{x}_{p \in P}\}$ with number of $c_{\mathrm{in}}$ channels, the position of each feature $p \in P$ sparsely distributed in the 3D spatial space. We process these features by a convolution with kernel weights $\mathrm{w} \in \mathbb{R}^{K \times c_{\mathrm{in}} \times c_{\mathrm{out}}}$. For example, in the 3D coordinate space, w contains $c_{\mathrm{in}} \times c_{\mathrm{out}}$ spatial kernels with size 3 and $|K| = 3^3$. The convolution process at the output position $\bar{p}$ is represented as

$$\mathrm{y}_{\bar{p}} = \sum_{k \in K} \mathrm{w}_k \cdot \mathrm{x}_{\bar{p}+k}, \tag{1}$$

where $k$ is an 3 dimensional offset distance from $\bar{p}$. $\bar{p} + k$ is the corresponding location around center $\bar{p}$. It enumerates all discrete locations in the kernel space $K$ where sparse features $\mathrm{x}_{\bar{p}+k}$ exist.

**Convolution types.** Regular sparse convolutions enlarge their output positions by dilating each input position $p \in P$ to the kernel shape $K$. It dramatically decreases feature density and increases the computational burden. Thus, regular sparse convolutions with stride 2 are only employed at the first layer in each stage for down-sampling.

Submanifold sparse convolutions, in contrast, fix their output positions $\bar{p}$ identical to input positions $\bar{p} \in P$, which preserves the computation cost at low levels. Due to the efficiency advantage, submanifold sparse convolutions dominate most layers of 3D CNNs, except for the down-sampling layers. However, limited by the small local scope, it misses sufficient information flow, for the spatially disconnected features, as shown in Fig. 1 (left). Naturally, increasing kernel sizes as in Fig. 1 (middle) becomes a potential solution to relieve this issue.

**Obstacles in 3D large kernels.** *Efficiency* is the first issue in 3D large-kernel CNNs. When we increase kernel sizes, the amount of parameters and computational burden have more rapid growth ratio than those of 2D CNNs. For instance, the parameter amount in one 3D convolutional layer increases from $27 \cdot C_{in \cdot out}$ ($3^3$) to $343 \cdot c_{in} \cdot c_{out}$ ($7^3$), as its kernel size from 3 to 7. It suggested that naively enlarging 3D kernels is unreasonable.

*Optimization difficulty* is the second obstacle. The booming model size requires sufficient data amount for learning. However, 3D datasets are commonly not so large-scale as 2D benchmarks. For example, Microsoft COCO [34] and ImageNet [12] contain hundreds of thousands or millions of images, while 3D datasets commonly contain only thousands of scenes. Moreover, due to the sparsity of input, not all weights in 3D CNNs get an optimization chance in each training step, as shown in Fig. 1. This problem is more serious in large models where more weights are missed on average. Tab. 2 reveals this issue when increasing kernel sizes from $3^3$ to $5^3$ in MinkowskiNet-34 [10].

*Previous experienced knowledge* are not helpful in 3D large-kernel CNNs. Recent 2D literature [15, 38] presents some components that are necessary or helpful in large-kernel CNNs, including depth-wise convolutions [26], layer normalization [1], GELU [25] activation. We examine these popularly used components and find them ineffective or even harmful in 3D CNNs. To demonstrate this,
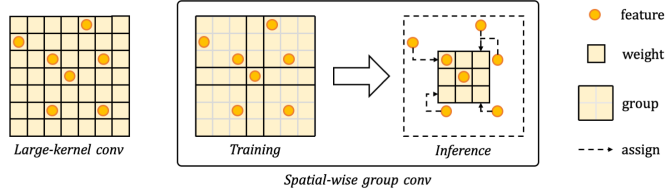
4

Figure 4: Illustration on spatial-wise group convolution. In the plain large-kernel convolution, due to the sparsity of input features, only a small proportion of weights are involved and updated in each training step. The spatial-wise group convolution relieves this issue by sharing weights in among spatial neighbors, increasing the chance of weights being optimized. During inference, we reformulate the spatial-wise convolutions into small ones, with feature assignment in a large scope.

we use MinkowskiNet-14 [10] to benchmark and train on the 3D semantic segmentation dataset, ScanNetv2 [11]. We substitute the original activation function from ReLU to GELU, the batch normalization to layer normalization, and the plain convolution to depth-wise, respectively. Tab. 1 shows that directly replacing these components results in performance drops, no matter on small or large kernel networks. 3D large-kernel CNNs require more specific designs.

## 3.2 Designs of 3D Large-kernel CNNs

Naively applying large-kernel convolutions to 3D CNNs usually harms the performance and efficiency with no benefits. In this section, we summarize our solutions for effective large kernels in 3D CNNs.

**Spatial-wise group convolution.** A standard convolution kernel can be viewed as a 3-dimensional matrix, which consists of input channels $C_{\text{in}}$, output channels $C_{\text{out}}$, and the spatial kernel dimension $K$. Taking kernel size as $k$, the volume of spatial kernel dimension $K$ equals $k^2$ for 2D convolutions and $k^3$ for 3D convolutions. Depth-wise convolutions share weights along channel dimensions, whose group numbers are equal to input channels. Point-wise convolutions fix kernel size as 1, which is commonly used as a conjunction layer to depth-wise convolutions to adjust output channels. We illustrate the kernel shapes of different convolution types in Fig. 3.

Different from 2D literature, which relies on depth-wise convolutions for performance boosting [22, 38] and accuracy-FLOPs trade-offs [15], we empirically find that depth-wise convolutions are non-beneficial in 3D tasks (shown in Tab. 1), no matter whether point-wise convolutions are included. Specially, we propose another kind of group convolution, spatial-wise convolution, for 3D large-kernel CNNs. It shares weights among spatial dimension $K$ on convolutional kernels, instead of channel dimension (also different from SGC [66], which partitions spatial groups on input features.). Specially, in Eq. 1, $w_k$ shares the same values in a local area $k \in K$. For example, we can group the original large-kernel convolution, from $7 \times 3$, into $3 \times 3$, by sharing weights among spatial neighbors. Spatial-wise group convolution is specially designed for 3D tasks and has the following advantages in terms of efficiency and performance.

The *efficiency* of 3D sparse large-kernel convolution is a severe problem, as shown in Tab. 3. When we naively enlarge the kernel size from $3^3$ to $7^3$ on MinkowskiNet-14 [10], both model size and runtime have a sharp increase by multiple times. In contrast, the spatial-wise group convolution avoids the increase in parameters and introduces limited runtime overhead. The advantage of our runtime over plain large kernels attributes to the implementation of sparse convolution based on the *gather-matmul-scatteradd* pipeline. In short, for a spatial-wise group convolution, we can directly use the small-kernel layer during inference, but enlarge its feature-assign areas to the large-kernel scope. In other words, thanks to the weight-sharing operation, it largely saves the *scatteradd* operations, *e.g.*, from $7^3$ to $3^3$ times. We integrated this implementation upon the open-sourced *spconv* library and omit the details here. Tab. 3 shows that our implementation on spatial-wise group convolution is far more efficient than the plain large-kernel baseline.

The *optimization* difficulty of large-kernel CNNs is relieved with the help of spatial-wise group convolutions. It constrains model sizes to a limited number, which correspondingly avoids requiring a large amount of data. Moreover, 3D input data or features, *e.g.*, voxels, are sparse and non-uniformly distributed, different from 2D images, which are dense and structured. This intrigues a situation that
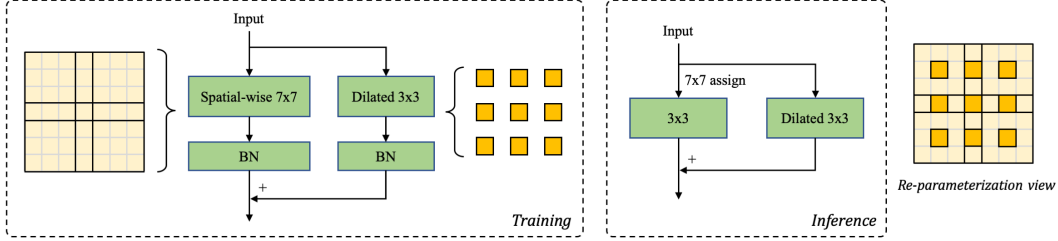
Figure 5: Structure of spatial-wise large-kernel block (SW-LK block). It consists of a large-kernel spatial-wise group convolution and a parallel branch. The convolutional parallel branch is introduced to make up the detail-capturing ability of large kernels. Dilation is required to put its kernel positions onto the group centers, in the view of re-parameterization.



(a) SW-LKNet for 3D semantic segmentation
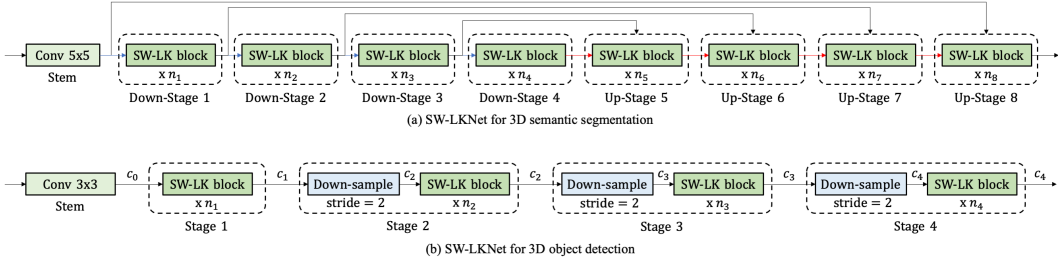


(b) SW-LKNet for 3D object detection

Figure 6: Architectures of spatial-wise large-kernel networks (LargeKernel3D). The network for 3D semantic segmentation follows the U-Net structure. We substitute all convolutional blocks, except for the down-sampling and up-sampling layers, to our SW-LK blocks in MinkowskiNet [10]. The network for 3D object detection consists of 1 stem layer and 4 stages, where we also replace the common blocks to be our SW-LK blocks. More details are provided in the *appendix*.

not all weights are involved and updated in each training step, as shown in Fig. 4. And the proportion of unused weights increases as kernels become larger, degrading the optimization of 3D large-kernel CNNs. When we include spatial-wise groups into large-kernel layers by weight sharing, weights have more chance to be optimized during training. In other words, grouped weights will all be updated, even if only one feature falls in the local group. Thus, spatial-wise group convolutions outperform its plain large-kernel counterparts (Tab. 3), especially on large models (Tab. 2).

**A parallel convolutional branch.** Although spatial-wise group convolutions provide an efficient alternative for large kernels, it potentially blurs detailed information. Because weight-sharing inevitably involves feature pooling. To relieve this issue, we introduce a parallel convolutional branch to capture details. Dilation might be required to make kernels placed inside the spatial groups. As shown in Fig. 5, upon the spatial-wise $7\times7$ kernel, a $3\times3$ convolutional layer with dilation 2 serves as a parallel branch. Batch Normalization (BN) layers are added in each branch. With the help of this branch, the originally blurred local areas are split in the re-parameterization [14] view.

During inference, we only fuse BN layers and keep two parallel branches separate [15]. As shown in Fig. 5, we convert the original spatial-wise $7\times7$ layer into $3\times3$ with features assigned in $7\times7$ kernel scopes. This small-kernel parallel branch introduces very limited additional computations, as shown in Tab. 3. These two paralleled small-kernel convolutions are more efficient than a large one.

### 3.3 Large-kernel Architecture

Following the above designs and observations, we introduce the architectures of our large-kernel 3D CNNs. We readily replace the plain 3D convolution blocks with our spatial-wise large-kernel (SW-LK) blocks in existing 3D backbones for semantic segmentation and object detection in Fig. 6.

**3D semantic segmentation.** The architectures of LargeKernel3D for 3D semantic segmentation resemble MinkowskiNet [10], including 1 stem layer and 8 stages. The stem layer is a $5\times5$ submanifold

6

Table 1: Ablations on using 2D techniques in 3D CNNs on ScanNetv2. Each is used separately.

| Method | Kernel | mIoU |
|---|---|---|
| MinkowskiNet-14 | | 67.0 |
| + GELU | | 66.3$_\downarrow$ |
| + layer norm | $3^3$ | 65.0$_\downarrow$ |
| + *depth-wise* conv | | 46.8$_\downarrow$ |
| MinkowskiNet-14 (K7) | | 68.6 |
| + GELU | | 67.8$_\downarrow$ |
| + layer norm | $7^3$ | 66.1$_\downarrow$ |
| + dilated conv $3^3$ | | 62.3$_\downarrow$ |
| + *depth-wise* conv | | 56.4$_\downarrow$ |

Table 2: Ablations on large models on ScanNetv2. $\downarrow$ ($\uparrow$) means accuracy increase (decrease).

| Method | Kernel | mIoU |
|---|---|---|
| MinkowskiNet-34 | $3^3$ | 71.7 |
| | $5^3$ | 70.7$_\downarrow$ |
| | $7^3$ | 68.6$_\downarrow$ |
| MinkowskiNet-34 + *depth-wise* conv | $3^3$ | 70.6$_\downarrow$ |
| | $5^3$ | 70.0$_\downarrow$ |
| | $7^3$ | 68.7$_\updownarrow$ |
| LargeKernel3D-34 | $5^3 \to 3^3$ | 72.2$_\uparrow$ |
| | $7^3 \to 3^3$ | 73.2$_\uparrow$ |

convolution. The first 4 stages account for down-sampling and the last 4 stages are for up-sampling. As for the U-Net structure, the first 4 features are concatenated into the up-sampling stages. We replace the original layers in MinkowskiNets [10] by $7 \times 7$ SW-LK blocks, where spatial-wise $7 \times 7$ convolutions is grouped into $3\times3$ splits. In addition to kernel sizes, model capacities, *e.g.*, MinkowskiNet-14A and MinkowskiNet-34C, are determined by channels $\{c_0, c_1, c_2, c_3, c_4, c_5, c_6, c_7, c_8\}$ and block numbers $\{n_1, n_2, n_3, n_4, n_5, n_6, n_7, n_8\}$.

**3D object detection.** The typical backbone networks in 3D object detectors [52, 13, 63] consist of 1 stem layer and 4 stages. A $3\times3$ submanifold convolution [21] layer serves as the stem. In each stage, except the first one, there is a sparse convolutional layer with stride 2 for down-sampling. We substitute other plain blocks with SW-LK blocks. In addition to the kernel sizes $\{k_1, k_2, k_3, k_4\}$, the overall architectures are determined by channels $\{c_0, c_1, c_2, c_3, c_4\}$ and block numbers $\{n_1, n_2, n_3, n_4\}$. We keep these hyper-parameters unchanged to the default settings in baseline detectors [13, 63]. We include these detailed numbers in the *appendix*.

## 4  Experiments

We conduct ablations and comparisons on 3D semantic segmentation and object detection. Additional ablations, *e.g.*, group manners and kernel size selections are included in the *appendix*.

### 4.1  Setups and Implementations

**ScanNetv2 for 3D semantic segmentation.** ScanNetv2 [11] is a large-scale benchmark in 3D semantic segmentation. It contains 1201 indoor RGB-D scenes for training, 312 scenes for validation, and 100 scenes for testing. Semantic labels in 20 categories are annotated. We train our models with the default settings in MinkowskiNet [10], including training hyper-parameters, and data augmentations. We evaluate models in the mean Intersection over Union (mIoU) metric. For ablation studies, we apply a tiny version of MinkowskiNet-14A [10] network. Detailed experimental settings, including network structures and training details, are listed in the *appendix*.

**KITTI and nuScenes for 3D object detection.** KITTI and nuScenes datasets are both popular 3D object detection benchmarks. KITTI dataset [19] contains 7,481 samples and 7,518 testing samples. The training set is typically split into a *train* set with 3,717 samples and a *val* set with 3,769 samples. The mean average precision metric for 3D bounding box ($AP_{3D}$) results is used for evaluation. It is commonly calculated with recall 40 positions (recall 40) or 11 positions (recall 11). nuScenes is a larger dataset and contains 1,000 driving sequences in total. Among them, there are 700 scenes split for training, 150 scenes for validation, and 150 scenes for testing. It contains LIDAR, camera, and radar sources with a complete $360^\circ$ environment. We evaluate our methods on both LIDAR-only and LIDAR-RGB fusion settings. The main evaluation metrics are mAP and nuScenes detection score (NDS). In experiments, we validate our networks on state-of-the-art detection frameworks of Voxel R-CNN [13] on KITTI [19], and CenterPoint [63] on nuScenes [3]. Other settings follow baselines by default [13, 63], including training and network hyper-parameters.

Table 3: Ablations on spatial-wise group convolutions and the parallel branch.

| Method | Kernel | # Params | Runtime | mIoU |
|---|---|---|---|---|
| MinkowskiNet-14 | $3^3$ | 1.87M | 45ms | 67.0 |
| MinkowskiNet-14 (K7) | $7^3$ | 21.21M | 220ms | 68.6 |
| + *spatial-wise* group conv | $7^3 \rightarrow 3^3$ | 1.87M | 72ms | 69.3 |
| + parallel branch (Ours) | | 2.77M | 97ms | **69.7** |

Table 4: Comparisons on ScanNetv2 mIoU on 3D semantic segmentation. [†] Sliding-window testing.

| Method | val | *test* |
|---|---|---|
| PointCNN [31] | - | 45.8 |
| PointNet++ [47] | 53.5 | 55.7 |
| RandLA-Net [27] | - | 64.5 |
| PointConv [58] | 61.0 | 66.6 |
| PointASNL [59] | 63.5 | 66.6 |
| KPConv [54] | 69.2 | 68.6 |
| FusionNet [65] | - | 68.8 |
| Point Transformer[†] [67] | 70.6 | - |
| Fast Point Transformer [43] | 72.1 | - |
| SparseConvNet [20] | 69.3 | 72.5 |
| MinkowskiNet-42 [10] | - | 73.4 |
| Stratified Transformer[†] [29] | 74.3 | 73.7 |
| MinkowskiNet-34 (baseline) | 71.7 | - |
| LargeKernel3D | 73.2 | **73.9** |

Table 5: Comparison on KITTI *val* split in $AP_{3D}$ in Recall 11 for the *Car* category.

| Method | Easy | **Mod.** | Hard |
|---|---|---|---|
| VoxelNet [45] | 81.97 | 65.46 | 62.85 |
| PointPillars [30] | 86.62 | 76.06 | 68.91 |
| SECOND [60] | 88.61 | 78.62 | 77.22 |
| Point R-CNN [51] | 88.88 | 78.63 | 77.38 |
| Part-$A^2$ [53] | 89.47 | 79.47 | 78.54 |
| 3DSSD [62] | 89.71 | 79.45 | 78.67 |
| Pointformer [42] | 90.05 | 79.65 | 78.89 |
| SA-SSD [23] | 90.15 | 79.91 | 78.78 |
| PV-RCNN [52] | 89.35 | 83.69 | 78.70 |
| VoTr-TSD [41] | 89.04 | 84.04 | 78.68 |
| Pyramid-PV [40] | 89.37 | 84.38 | 78.84 |
| Focals Conv [9] | 89.52 | 84.93 | 79.18 |
| Voxel R-CNN [13] | 89.41 | 84.52 | 78.93 |
| LargeKernel3D | 89.52 | **85.07** | 79.32 |

## 4.2 Ablation Studies

**Effectiveness of techniques in 2D CNNs.** We first validate the techniques that proven effective on 2D CNNs, including layer normalization [1], depth-wise convolution [26, 38], and GeLU [25]. We conduct experiments on a tiny MinkowskiNet-14 [10] on the ScanNet [11] semantic segmentation dataset, as shown in Tab. 1. MinkowskiNet-14 (K7) is a modified network to MinkowskiNet-14, with kernel sizes in all stages enlarged to 7. All these techniques bring no improvements to baseline networks, especially depth-wise convolution, which severely harms the performance by a large margin. In addition, we also perform depth-wise convolutions on large models, MinkowskiNet-34 [10], in Tab. 2. It shows that depth-wise convolutions are also non-beneficial on large models.

Another alternative for enlarging kernel sizes is dilated convolution [4], which is a common component in 2D semantic segmentation. However, Tab. 1 shows dilated convolution (from kernel $3^3$ to $7^3$) in 3D semantic segmentation largely degrades the performance, compared to the plain convolution baseline. We suppose that dilation factorizes the original dense kernels into sparse ones. While input 3D features are also sparse, uncovered information will be missed during feature aggregation.

**Spatial-wise group convolution and parallel branch.** We validate our designs of spatial-wise group convolution and parallel branch upon the tiny MinkowskiNet-14 [10] in Tab. 3. We evaluate our method on the ScanNet 3D semantic segmentation dataset in both accuracy and efficiency. Runtimes are measured on an NVIDIA 2080Ti GPU. Naively enlarging the kernel size of MinkowskiNet-14 [10] to kernel 7 indeed improves the performance by 1.6%. However, it dramatically increases computational burdens. In addition, naively using large-kernel in large model MinkowskiNet-34 [10] is shown harmful as in Tab. 2. In contrast, ours present better accuracy by 2.7% mIoU with linear parameters and runtime overhead, rather than the cubic increase of naive large kernels.

**Ablations on large models.** In addition to the above small-model experiments, we also conduct ablations on the large model MinkowskiNet-34 [10] in Tab. 2. All models are trained with the same hyper-parameters in ablation studies. By naively enlarging kernel size in MinkowskiNet-34, the performance drops by a large margin. Depth-wise convolution also brings no benefit. In contrast, our LargeKernel3D-34 gradually improves the performance from 71.7% to 73.2% mIoU.

Table 6: Comparison with other methods on nuScenes *test* split.

| Method | NDS | mAP | Car | Truck | Bus | Trailer | C.V. | Ped | Mot | Byc | T.C. | Bar |
|--------|-----|-----|-----|-------|-----|---------|------|-----|-----|-----|------|-----|
| PointPillars [30] | 45.3 | 30.5 | 68.4 | 23.0 | 28.2 | 23.4 | 4.1 | 59.7 | 27.4 | 1.1 | 30.8 | 38.9 |
| 3DSSD [62] | 56.4 | 42.6 | 81.2 | 47.2 | 61.4 | 30.5 | 12.6 | 70.2 | 36.0 | 8.6 | 31.1 | 47.9 |
| CBGS [68] | 63.3 | 52.8 | 81.1 | 48.5 | 54.9 | 42.9 | 10.5 | 80.1 | 51.5 | 22.3 | 70.9 | 65.7 |
| CenterPoint [63] | 65.5 | 58.0 | 84.6 | 51.0 | 60.2 | 53.2 | 17.5 | 83.4 | 53.7 | 28.7 | 76.7 | 70.9 |
| HotSpotNet [6] | 66.0 | 59.3 | 83.1 | 50.9 | 56.4 | 53.3 | 23.0 | 81.3 | 63.5 | 36.6 | 73.0 | 71.6 |
| CVCNET [5] | 66.6 | 58.2 | 82.6 | 49.5 | 59.4 | 51.1 | 16.2 | 83.0 | 61.8 | 38.8 | 69.7 | 69.7 |
| TransFusion [2] | 70.2 | 65.5 | 86.2 | 56.7 | 66.3 | 58.8 | 28.2 | 86.1 | 68.3 | 44.2 | 82.0 | 78.2 |
| Focals Conv [9] | 70.0 | 63.8 | 86.7 | 56.3 | 67.7 | 59.5 | 23.8 | 87.5 | 64.5 | 36.3 | 81.4 | 74.1 |
| Focals Conv-F[‡] [9] | 73.6 | 70.1 | 87.5 | 60.0 | 69.9 | 64.0 | 32.6 | 89.0 | 81.1 | 59.2 | 85.5 | 71.8 |
| LargeKernel3D | 70.5 | 65.3 | 85.9 | 55.3 | 66.2 | 60.2 | 26.8 | 85.6 | 72.5 | 46.6 | 80.0 | 74.3 |
| LargeKernel3D[‡] | 72.8 | 68.8 | 87.3 | 59.1 | 68.5 | 65.6 | 30.2 | 88.3 | 77.8 | 53.5 | 82.4 | 75.0 |
| LargeKernel3D-F[‡] | **74.2** | **71.1** | 88.1 | 60.3 | 69.1 | 66.5 | 34.3 | 89.6 | 82.0 | 60.3 | 85.7 | 75.5 |

[‡] Flipping and rotation testing-time augmentations.

### 4.3 3D Semantic Segmentation

We make comparisons with other 3D semantic segmentation methods on the ScanNetv2 dataset in Tab. 4. Our method surpasses others on *test* split. MinkowskiNet [10] is already a state-of-the-art method in ScanNetv2. Our SW-LK block further improves its performance. Our method is also superior to transformer-based methods [67, 43]. Only Stratified Transformer [29] presents better mIoU on *val* split than ours. However, it uses an unfair sliding-window testing technique[†]. It splits integrated scenes into overlapped parts to test each part one by one and then ensembles the results together, taking hours on evaluation. Moreover, our results are better on the *test* split.

### 4.4 3D Object Detection

**KITTI.** On the KITTI [19] dataset, we evaluate our model on the *val* split upon Voxel R-CNN [13]. We report the results in the recall 40 metric on *Car*, which is the main category in this dataset. Our results outperform other LIDAR-only methods, including the transformer-based method, VoTr-TSD [41], and the sparse convolutional network Focals Conv [9]. Voxel R-CNN is already a state-of-the-art model on KITTI. We improve this strong baseline to 85.07% $AP_{3D}$.

**nuScenes.** We compare our LargeKernel3D upon CenterPoint [63] with previous methods on the *test* split of the nuScenes [3] dataset in Tab. 6. LargeKernel3D improves CenterPoint [63] to 70.5% and 72.8% NDS, with and without test augmentations, both outperforming other LIDAR-only methods. The multi-modal modal LargeKernel3D-F further improves to 74.2% NDS. The results on the *val* split and the mutli-modal network setting are included in the *appendix*.

## 5 Conclusion

This paper studies large-kernel designs for 3D convolutional networks, which have essential differences from the solutions in 2D CNNs. We present spatial-wise group convolution that is specifically designed for 3D large kernels. It effectively resolves the efficiency and optimization issues in plain 3D large-kernel CNNs. Based on this design, we further propose the SW-LK block and the corresponding LargeKernel3D for 3D semantic segmentation and object detection. Our 3D large-kernel networks achieve non-trivial improvements on both semantic segmentation and object detection benchmarks. For the first time, we show that 3D large kernels are feasible and essential. We believe that these improvements come from the enlarged effective receptive field. We hope that our findings could advance the studies of 3D network design. Additional experiments are in the *appendix*.

**Limitations.** LargeKernel3D relies on the pre-defined spatial kernel sizes in 3D semantic segmentation and object detection benchmarks. This selected value might be sub-optimal for other datasets or tasks, *e.g.*, panoptic segmentation, depending on the overall scene sizes and data distribution. Other techniques, *e.g.*, neural architecture search (NAS) [69, 8], might be helpful for future work.

**Boarder impacts.** The proposed method serves backbone networks for various 3D tasks, which might include tasks or datasets involving negative societal impacts, while the network is innocent.

# References

[1] Lei Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *CoRR*, abs/1607.06450, 2016.

[2] Xuyang Bai, Zeyu Hu, Xinge Zhu, Qingqiu Huang, Yilun Chen, Hongbo Fu, and Chiew-Lan Tai. Transfusion: Robust lidar-camera fusion for 3d object detection with transformers. *CoRR*, abs/2203.11496, 2022.

[3] Holger Caesar, Varun Bankiti, Alex H. Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *CVPR*, pages 11618–11628, 2020.

[4] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *T-PAMI*, 40(4):834–848, 2017.

[5] Qi Chen, Lin Sun, Ernest Cheung, and Alan L. Yuille. Every view counts: Cross-view consistency in 3d object detection with hybrid-cylindrical-spherical voxelization. In *NeurIPS*, 2020.

[6] Qi Chen, Lin Sun, Zhixin Wang, Kui Jia, and Alan L. Yuille. Object as hotspots: An anchor-free 3d object detection approach via firing of hotspots. In *ECCV*, volume 12366, pages 68–84, 2020.

[7] Xuanyao Chen, Tianyuan Zhang, Yue Wang, Yilun Wang, and Hang Zhao. FUTR3D: A unified sensor fusion framework for 3d detection. *CoRR*, abs/2203.10642, 2022.

[8] Yukang Chen, Tong Yang, Xiangyu Zhang, Gaofeng Meng, Xinyu Xiao, and Jian Sun. Detnas: Backbone search for object detection. In *NeurIPS*, pages 6638–6648, 2019.

[9] Yukang Chen, Yanwei Li, Xiangyu Zhang, Jian Sun, and Jiaya Jia. Focal sparse convolutional networks for 3d object detection. *CoRR*, abs/2204.12463, 2022.

[10] Christopher B. Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *CVPR*, pages 3075–3084, 2019.

[11] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *CVPR*, 2017.

[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255, 2009.

[13] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel R-CNN: towards high performance voxel-based 3d object detection. In *AAAI*, pages 1201–1209, 2021.

[14] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: Making vgg-style convnets great again. In *CVPR*, pages 13733–13742, 2021.

[15] Xiaohan Ding, Xiangyu Zhang, Yizhuang Zhou, Jungong Han, Guiguang Ding, and Jian Sun. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. *CoRR*, abs/2203.06717, 2022.

[16] Xiaoyi Dong, Jianmin Bao, Dongdong Chen, Weiming Zhang, Nenghai Yu, Lu Yuan, Dong Chen, and Baining Guo. Cswin transformer: A general vision transformer backbone with cross-shaped windows. *CoRR*, abs/2107.00652, 2021.

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021.

[18] Lue Fan, Ziqi Pang, Tianyuan Zhang, Yu-Xiong Wang, Hang Zhao, Feng Wang, Naiyan Wang, and Zhaoxiang Zhang. Embracing single stride 3d object detector with sparse transformer. *CoRR*, abs/2112.06375, 2021.

[19] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The KITTI dataset. *Int. J. Robotics Res.*, 32(11):1231–1237, 2013.

[20] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, pages 9224–9232, 2018.

[21] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *CVPR*, pages 9224–9232, 2018.

[22] Qi Han, Zejia Fan, Qi Dai, Lei Sun, Ming-Ming Cheng, Jiaying Liu, and Jingdong Wang. On the connection between local attention and dynamic depth-wise convolution. In *ICLR*, 2021.

[23] Chenhang He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *CVPR*, pages 11870–11879, 2020.

[24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016.

[25] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *CoRR*, abs/1606.08415, 2016.

[26] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *CoRR*, abs/1704.04861, 2017.

[27] 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 *CVPR*, pages 11105–11114, 2020.

[28] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, volume 37, pages 448–456, 2015.

[29] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. *CoRR*, abs/2203.14508, 2022.

[30] Alex H. Lang, Sourabh Vora, Holger Caesar, Lubing Zhou, Jiong Yang, and Oscar Beijbom. Pointpillars: Fast encoders for object detection from point clouds. In *CVPR*, pages 12697–12705, 2019.

[31] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. In *NeurIPS*, pages 828–838, 2018.

[32] Yanwei Li, Xiaojuan Qi, Yukang Chen, Liwei Wang, Zeming Li, Jian Sun, and Jiaya Jia. Voxel field fusion for 3d object detection. In *CVPR*, 2022.

[33] Tingting Liang, Hongwei Xie, Kaicheng Yu, Zhongyu Xia, Zhiwei Lin, Yongtao Wang, Tao Tang, Bing Wang, and Zhi Tang. Bevfusion: A simple and robust lidar-camera fusion framework. *CoRR*, abs/2205.13790, 2022.

[34] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft COCO: common objects in context. In *ECCV*, volume 8693, pages 740–755, 2014.

[35] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. Swin transformer V2: scaling up capacity and resolution. *CoRR*, abs/2111.09883, 2021.

[36] 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 *ICCV*, pages 9992–10002, 2021.

[37] Zhijian Liu, Haotian Tang, Alexander Amini, Xinyu Yang, Huizi Mao, Daniela Rus, and Song Han. Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation. *CoRR*, abs/2205.13542, 2022.

[38] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *CoRR*, abs/2201.03545, 2022.

[39] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard S. Zemel. Understanding the effective receptive field in deep convolutional neural networks. In *NeurIPS*, pages 4898–4906, 2016.

[40] Jiageng Mao, Minzhe Niu, Haoyue Bai, Xiaodan Liang, Hang Xu, and Chunjing Xu. Pyramid R-CNN: towards better performance and adaptability for 3d object detection. In *ICCV*, 2021.

[41] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *ICCV*, 2021.

[42] Xuran Pan, Zhuofan Xia, Shiji Song, Li Erran Li, and Gao Huang. 3d object detection with pointformer. In *CVPR*, pages 7463–7472, 2021.

[43] Chunghyun Park, Yoonwoo Jeong, Minsu Cho, and Jaesik Park. Fast point transformer. *CoRR*, abs/2112.04702, 2021.

[44] Chao Peng, Xiangyu Zhang, Gang Yu, Guiming Luo, and Jian Sun. Large kernel matters - improve semantic segmentation by global convolutional network. In *CVPR*, pages 1743–1751, 2017.

[45] Charles R. Qi, Or Litany, Kaiming He, and Leonidas J. Guibas. Deep hough voting for 3d object detection in point clouds. In *ICCV*, pages 9276–9285, 2019.

[46] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 77–85, 2017.

[47] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *NeurIPS*, pages 5099–5108, 2017.

[48] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? In *NeurIPS*, pages 12116–12128, 2021.

[49] Yongming Rao, Wenliang Zhao, Zheng Zhu, Jiwen Lu, and Jie Zhou. Global filter networks for image classification. In *NeurIPS*, pages 980–993, 2021.

[50] David W. Romero, Anna Kuzina, Erik J. Bekkers, Jakub M. Tomczak, and Mark Hoogendoorn. Ckconv: Continuous kernel convolution for sequential data. *CoRR*, abs/2102.02611, 2021.

[51] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Pointrcnn: 3d object proposal generation and detection from point cloud. In *CVPR*, pages 770–779, 2019.

[52] 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 *CVPR*, pages 10526–10535, 2020.

[53] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From points to parts: 3d object detection from point cloud with part-aware and part-aggregation network. *T-PAMI*, 43(8):2647–2664, 2021.

[54] 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 *ICCV*, pages 6410–6419, 2019.

[55] Asher Trockman and J. Zico Kolter. Patches are all you need? *CoRR*, abs/2201.09792, 2022.

[56] Sourabh Vora, Alex H. Lang, Bassam Helou, and Oscar Beijbom. Pointpainting: Sequential fusion for 3d object detection. In *CVPR*, pages 4603–4611, 2020.

[57] Chunwei Wang, Chao Ma, Ming Zhu, and Xiaokang Yang. Pointaugmenting: Cross-modal augmentation for 3d object detection. In *CVPR*, pages 11794–11803, 2021.

[58] Wenxuan Wu, Zhongang Qi, and Fuxin Li. Pointconv: Deep convolutional networks on 3d point clouds. In *CVPR*, pages 9621–9630, 2019.

[59] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *CVPR*, pages 5588–5597, 2020.

[60] Yan Yan, Yuxing Mao, and Bo Li. SECOND: sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018.

[61] Jianwei Yang, Chunyuan Li, Pengchuan Zhang, Xiyang Dai, Bin Xiao, Lu Yuan, and Jian-feng Gao. Focal self-attention for local-global interactions in vision transformers. *CoRR*, abs/2107.00641, 2021.

[62] Zetong Yang, Yanan Sun, Shu Liu, and Jiaya Jia. 3dssd: Point-based 3d single stage object detector. In *CVPR*, pages 11037–11045, 2020.

[63] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-based 3d object detection and tracking. In *CVPR*, pages 11784–11793, 2021.

[64] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Multimodal virtual point 3d detection. In *NeurIPS*, 2021.

[65] Feihu Zhang, Jin Fang, Benjamin W. Wah, and Philip H. S. Torr. Deep fusionnet for point cloud semantic segmentation. In *ECCV*, volume 12369, pages 644–663, 2020.

[66] Jiahui Zhang, Hao Zhao, Anbang Yao, Yurong Chen, Li Zhang, and Hongen Liao. Efficient semantic scene completion network with spatial group convolution. In *ECCV*, pages 733–749, 2018.

[67] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip H. S. Torr, and Vladlen Koltun. Point transformer. In *ICCV*, pages 16239–16248, 2021.

[68] Benjin Zhu, Zhengkai Jiang, Xiangxin Zhou, Zeming Li, and Gang Yu. Class-balanced grouping and sampling for point cloud 3d object detection. *CoRR*, abs/1908.09492, 2019.

[69] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable architectures for scalable image recognition. In *CVPR*, pages 8697–8710, 2018.
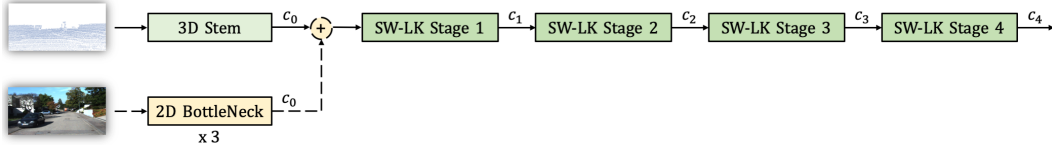
Figure S - 7. Architectures of LargeKernel3D with image fusion for 3D object detection.
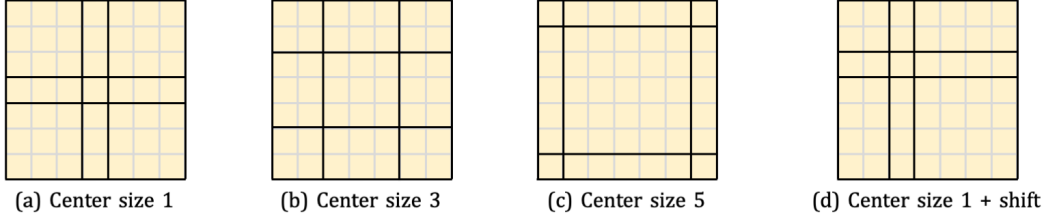


(a) Center size 1     (b) Center size 3     (c) Center size 5     (d) Center size 1 + shift

Figure S - 8. The kernel group manner choices. We study the center sizes for group splitting and whether to shift the centers.

# Appendix

## A    More Implementation Details

### A.1    Data processing.

**ScanNetv2.** We convert point clouds into voxels as input data for the ScanNetv2 dataset. The voxelization sizes are all 0.02m for all $X$, $Y$, $Z$ axes. In terms of data augmentations, we exactly follow our baseline method, MinkowskiNet [10]. Specially, input data is randomly dropped out with a ratio 0.2. For spatial augmentations, we also conduct random horizontal flipping on input. For intensity augmentations, we conduct auto-contrast, color translation, and color jittering.

**KITTI.** We convert point clouds into voxels as input data for 3D object detectors. Point clouds is clipped into [0, 70.4m] for $X$ axis, [-40m,40m] for $Y$ axis, and [-3, 1]m for $Z$ axis, on KITTI [19] for ranges. The input voxel size is set as (0.05m, 0.05m, 0.1m). The data augmentations include random flipping, global scaling, global rotation, and ground-truth (GT) sampling [60] for the KITTI dataset. Random flipping is applied along the $X$ axis. Global scaling is sampled from the [0.95, 1.05] ratio. Global rotation is performed around the $Z$ axis. Rotation angle is sampled from [-45°, 45°]. Ground-truth sampling is to copy objects from other training data, and paste them onto the current scene. It enriches data variance during training.

**nuScenes.** We clip point clouds into [-54m, 54m] for both $X$ and $Y$ axes, and [-5m, 3m] for the $Z$ axis, on nuScenes [3]. The voxel size is set as (0.075m, 0.075m, 0.2m). Data augmentations includes random flipping, global scaling, global rotation, GT sampling [60], and an additional translation on the nuScenes [3] dataset. Random flipping is performed in $X$ and $Y$ axes. Rotation angle is sampled in [-45°, 45°]. Global scaling is conducted in the [0.9, 1.1] ratio. Translation noise is conducted on all three axes from the ratio [0, 0.5]. GT sampling is also conducted on the nuScenes.

### A.2    Training settings.

**ScanNetv2.** For models trained on the ScanNetv2 dataset, we train the network for 600 epochs with batch size 16. The learning rate is initialized as 0.1 and decays with a poly scheduler. We adopt the SGD optimizer. The momentum is set as 0.9. The hyper-parameters directly follow our baseline [10].

**KITTI.** We train the network for 80 epochs and batch size 16, for Voxel R-CNN [13] model on KITTI. The learning rate is initialized as 0.01 and decreases in the cosine annealing strategy. Gradient norms are clipped by 10. We adopt the Adam optimizer, with weight decay 0.01 and momentum 0.9.

Table S - 7. Improvement over CenterPoint on the nuScenes *val* split.

| Method | #Params | Runtime | mAP | NDS |
|---|---|---|---|---|
| CenterPoint | 9.0M | 96ms | 59.0 | 66.4 |
| LargeKernel3D | 9.66M | 118ms | **60.5** | **67.6** |

Table S - 8. Ablations on different kernel sizes and voxel sizes.

| Voxel size | Kernel $3^3$ | $5^3 \rightarrow 3^3$ | $7^3 \rightarrow 3^3$ |
|---|---|---|---|
| 0.01 m | <span style="color:red">64.5</span> | **69.7** | **70.4** |
| 0.02 m | **67.0** | 68.5 | 69.7 |

**nuScenes.** We train CenterPoint [63] on the nuScenes datasets for 20 epochs with batch size 32. This network is trained by Adam. The learning rate is set as 1e-3 and decays in the cosine annealing strategy to 1e-4. The weight decay is set as 0.01. And the gradient norms are clipped by 35.

### A.3 Network architecture settings.

**ScanNetv2.** We use two networks, MinkowskiNet-14 and MinkowskiNet-34 [10], for the ScanNetv2 dataset in the paper. In MinkowskiNet-14, we set the channels of each stage, $\{c_1, c_2, c_3, c_4, c_5, c_6, c_7, c_8\}$, as $\{16, 32, 64, 128, 64, 64, 32, 32\}$. There is one block in each stage. In MinkowskiNet-34, we set the channel numbers as $\{32, 64, 128, 256, 256, 128, 96, 96\}$. The block numbers, $\{n_1, n_2, n_3, n_4, n_5, n_6, n_7, n_8\}$, are $\{2, 3, 4, 6, 2, 2, 2, 2\}$. The meanings of these notations are shown in Fig. 6. LargeKernel3D-14 and LargeKernel3D-34 directly follows MinkowskiNet-14 and MinkowskiNet-34 respectively for these settings. They substitute the plain sparse convolutional blocks to the proposed SW-LK blocks with spatial size $7^3$ and groups $3^3$ in all stages.

**KITTI.** In the Voxel R-CNN [13] network, the channels of stages in the backbone network, $\{c_0, c_1, c_2, c_3, c_4\}$, are $\{16, 16, 32, 64, 64\}$. The block numbers in each stage, $\{n_1, n_2, n_3, n_4\}$, are $\{1, 2, 2, 2\}$. By default, each convolutional layer is followed by a batch normalization layer [28], and a ReLU activation. LargeKernel3D substitutes the plain sparse convolutional blocks to the proposed SW-LK blocks for the first 3 stages, *i.e.*, stage 1, 2, 3.

**nuScenes.** The backbone network of CenterPoint [63] has more layers. The channels $\{c_0, c_1, c_2, c_3, c_4\}$, equal to $\{16, 16, 32, 64, 128\}$. The block numbers in these stages, $\{n_1, n_2, n_3, n_4\}$, are $\{2, 2, 2, 2\}$. Compared to the networks in KITTI, each block contains two convolutional layers, with a residual connection, except the stem. LargeKernel3D also substitutes the plain blocks to SW-LK blocks for stage 1, 2, 3, without the last stage. Because the last stage has heavy channel numbers.

We also present the multi-modal network with our large kernel backbone, *i.e.*, LargeKernel3D-F in Tab. 6. As shown in Fig. 7, we conduct a direct voxel-wise summation between LIDAR and RGB features. The RGB branch only contains a conv-bn-relu-pooling stem and 3 residual bottlenecks [24]. We intentionally make the RGB branch lightweight to fully demonstrate the capacity of our large-kernel LIDAR backbone. Other details about multi-modal fusion refer to Focals Conv [9].

## B  Additional Experimental Results

### B.1  Ablations on kernel size selection.

We ablate the kernel size selection in Tab. 9. *Spatial* means the absolute spatial kernel size of our SW-LK blocks. *Groups* means the number of groups to split kernels for weight sharing. We conduct the experiment on LargeKernel3D-14 and the ScanNetv2 dataset, with the same training hyper-parameters as the main experiments. We empirically find that using the spatial kernel size $7^3$ is a good choice. Both $3^3$ and $5^3$ groups obtain satisfactory results. We choose the simple one, $3^3$ groups, as a default setting in the paper. Although fixing spatial size as $9^3$ with groups $7^3$ achieves similar accuracy, it is more complicated. Moreover, in the spatial size $9^3$, the performance decays as the group number decreases. We assume that the reason behind is weight sharing among too many parameters leads to severe feature blurring.

Table S - 9. Ablations on spatial size and groups.

| Spatial | Groups | mIoU |
|---------|--------|------|
| $3^3$ | $3^3$ | 67.0 |
| $5^3$ | $3^3$ | 68.5 |
| $7^3$ | $5^3$ | **69.7** |
| | $3^3$ | **69.7** |
| $9^3$ | $3^3$ | 68.9 |
| | $5^3$ | 69.3 |
| | $7^3$ | 69.6 |

Table S - 10. Ablations on group manners of spatial-wise group convolution on ScanNetv2.

| Centers | $1^3$ | $3^3$ | $5^3$ | $1^3$ + shift |
|---------|-------|-------|-------|---------------|
| mIoU | 69.3 | 67.4 | 65.0 | 63.7 |

Table S - 11. Improvements over CenterPoint on the nuScenes *val* split.

| Method | NDS | mAP | Car | Truck | Bus | Trailer | C.V. | Ped | Mot | Byc | T.C. | Bar |
|--------|-----|-----|-----|-------|-----|---------|------|-----|-----|-----|------|-----|
| CenterPoint [63] | 66.4 | 59.0 | 85.6 | 57.2 | 71.2 | 37.3 | 16.2 | 85.1 | 58.4 | 41.0 | 69.2 | 68.2 |
| LargeKernel3D | 67.6 | 60.5 | 85.8 | 59.0 | 72.8 | 40.0 | 19.2 | 85.6 | 61.3 | 43.6 | 70.4 | 67.5 |
| LargeKernel3D* | **69.1** | **63.3** | 85.9 | 59.8 | 73.9 | 42.5 | 22.5 | 85.4 | 68.5 | 56.1 | 71.5 | 66.7 |

* remove gt-sampling in the last 2 epochs.

## B.2 Relations between kernel sizes and voxel sizes.

In the voxel-based 3D semantic segmentation methods [10, 20], voxel size is usually set as {0.02m, 0.02m, 0.02m} in {$x$, $y$, $z$} axes as a default setting. In Tab. 8, we study the relations between kernel sizes and voxel sizes on our LargeKernel3D-14. Note that, Kernel 3 means the baseline model MinkowskiNet-14 [10]. Tt achieves an unsatisfied 64.5% mIoU, on the voxel size 0.01m. However, the performance is increased to 70.4% by 5.9% with our large-kernel design, which also outperforms its voxel-size 0.02 m counterpart. Small voxel sizes provide high-resolution representation for the 3D scenes. We show that it has more potential with the network with enough receptive field.

## B.3 Ablations on spatial group partition.

We also ablate the manners of group partition. As shown in Fig. 8, we study the center sizes and center shifting. This ablation study is also conducted on the LargeKernel3D-14 and the ScanNetv2 dataset. As shown in Tab. 10, we find that large center sizes and center shifting hurt the performance. As the center sizes enlarge from $1^3$ to $5^3$, the accuracy sharply drops. We suppose that large center sizes will blur individual features as weight sharing. In addition, center shifting is also harmful, because the asymmetric group manners miss the important central position. Therefore, we keep the center size as $1^3$, without shifting, as a default setting.

## B.4 Model complexity on nuScenes 3D object detection.

Tab. 7 presents the model complexity of our method upon the CenterPoint [63] baseline. It shows that the improvements on mAP are non-trivial from 59.0% to 60.5% mAP in a fair comparison, while the overhead of both parameters and runtime are limited.

## B.5 Further improvements on the nuScenes dataset.

We present further improvements on the nuScenes dataset by additional techniques in Tab. 11 and Tab. 6. These techniques are removing gt-sampling in the last training epochs and training with *val* split data for *test* server submission. These tricks have been used by some previous state-of-the-art methods [57, 9] for performance boosting, but not by all methods. Thus, we did not include these results in the main paper for a fair comparison. As shown in Tab. 6 and Tab. 11, LargeKernel3D achieves 63.3% mAP on the *val* split and 65.3% mAP on the *test split*, which are both very competitive results in this area. Note that other test augmentations might lead to additional improvements, but we think that these are unrelated to our main contributions.

Figure S - 9. Additional illustrations on effective receptive fields. Left - original images, mid - plain 3D CNNs, right - LargeKernel3D. It is best viewed in color and by zooming in.

## C  Visualizations

We provide additional visual comparisons between the plain 3D network and our LargeKernel3D in Fig. 9. It shares the same setting as the Fig. 2. These visualizations are implemented upon the Voxel R-CNN [13] network on the KITTI [19] dataset. In each group, the left one is the original image, the middle one is the ERFs of plain 3D CNNs, and the right one is the ERFs of our LargeKernel3D. We follow the definition of ERFs [39]. We calculate the gradient of every input voxel data regarding the feature of interest. It illustrates the intensity of feature changes as the input value changes. We normalize the gradient norms to [0, 1] and project these values onto the image plane with the calibration matrices.