

1 We thank all reviewers for their comments. All responses and changes will be  
 2 incorporated into the revision. Code will be released for full reproducibility.

3 **Baselines (R1, R3, R4).** We compared PVCNN with SSCN (sparse-conv)  
 4 and O-CNN (octree) as suggested by three reviewers. PVCNN is  $4.8\times$  faster  
 5 than SSCN with superior performance (Table 1). For O-CNN, the inference  
 6 is fast; however, it is more than  $10\times$  slower to preprocess the input data (*i.e.*,

7 condense the point cloud and construct the octree) and postprocess the output  
 8 (*i.e.*, refine the boundaries using the dense CRF). Also, as pointed out by R1, these two approaches are orthogonal to our  
 9 method and can be combined together: replacing the volumetric convolution with the sparse/octree-based convolution.

10 **Scale Normalization (R4).** We *only* normalize the point cloud in the voxel branch  
 11 while keeping the coordinates in the point branch unchanged; therefore, there is no  
 12 information loss after two branches are merged together. Without instance scale normal-  
 13 ization, the voxel grids are more than  $10\times$  sparser on average, and the volumetric  
 14 convolution is no longer effective to extract features. An alternative is to drop outlier  
 15 points that do not lie in any voxel grids, which will inevitably induce information loss  
 16 (see Table 2). Thus, the instance scale normalization is critical in the voxel branch.

17 **Devoxelization (R1, R4).** We compared different implementations of devoxeliza-  
 18 tion. From Table 2, the trilinear interpolation (w/ TI) performs better than the nearest  
 19 neighbor (w/o TI), which is because the points near the voxel boundaries will introduce  
 20 larger fluctuations to the gradient, making it harder to optimize, as mentioned by R4.

21 **Ablation Studies (R4).** More analyses are provided in Table 2, including different  
 22 feature fusion methods, voxel resolutions (also see Table 2 in the paper), and number  
 23 of convolutions per block. Our design choice is the best. We’ll add them to the paper.

24 **Evaluation on S3DIS (R4).** We follow exactly the same data processing and evalu-  
 25 ation protocol as PointCNN (L252) to make sure that the improvements are entirely from our proposed PVConv rather  
 26 than different evaluation protocols. We thank R4 for suggesting the ScanNet2 dataset which we will experiment on.

27 **Evaluation on KITTI (R1, R3).** We choose F-PointNet as our baseline, which is a popular open-source 3D detection  
 28 model. We do not compare with other state-of-the-art (voxel/point-based) models because the *region proposal network*  
 29 has a large impact on the final results; we want to remove the influence of different region proposal networks and only  
 30 evaluate the effectiveness of our *convolution primitive*. Besides, we do not present results on the testing set because  
 31 F-PointNet only provides the 2D region proposals on the training and validation set.

32 **Feature Visualization (R3).** We illustrate the voxel and point branch features from  
 33 the final PVConv in Figure 1. Note that warmer color represents larger magnitude.  
 34 It is interesting to see that the voxel branch captures large, continuous parts while the  
 35 point branch captures isolated, discontinuous details (*e.g.*, table legs, lamp necks).  
 36 The two branches provide complementary information and can be explained by the  
 37 fact that the convolution operation extracts features with continuity and locality.

38 **Statistical Significance (R1).** On KITTI, our model has already been evaluated for  
 39 **20 times** to reduce the variance (L266). On ShapeNet and S3DIS, our results have  
 40 relatively small variances: PVCNN achieves  $(86.13 \pm 0.04)\%$  in mIoU on ShapeNet,  
 41 and PVCNN++ achieves  $(58.95 \pm 0.08)\%$  in mIoU on S3DIS (trained for **4 times**  
 42 on both datasets). We want to emphasize that our goal is to achieve high accuracy *as well as* better efficiency. Some  
 43 improvements over PointCNN might be small in accuracy (0.1% in Table 1); however, the speedup is significant ( $2.7\times$ ).

44 **Improvement Breakdown (R1).** The speedup comes from our better algorithm with lower complexity rather than  
 45 the implementation: PVCNN can *theoretically* save the number of incontinentous memory accesses by  $k$  times (where  $k$   
 46 is the number of neighbors) and achieve better locality (L190-194). Both baselines and our models are implemented  
 47 using well-optimized deep learning libraries (PyTorch and TensorFlow). We do nothing further to optimize the speed.

48 **Sublinear Memory Growth (R1).** Asymptotically, the memory indeed grows cubically; however, the voxel resolution  
 49 of PVCNN is kept very low thanks to the high-resolution point branch (L246-247). With low resolutions (L242-244),  
 50 the memory consumption is dominated by the point branch, not the voxel branch. Thus, the memory grows sub-linearly.

51 **Paper Presentation (R1, R4).** We will revise our  
 52 paper to make it concise and objective. We would  
 53 like to clarify that the groups are assigned by the accu-  
 54 racy in Table 1 and 3 in the paper, and we put our  
 55 PVCNN’s into groups to highlight the speedup with  
 56 similar accuracy. As suggested, we also compared  
 57 *both* point-based and voxel-based models (see right)  
 58 to make Figure 4 in the paper more comprehensive.

	PVCNN	SSCN	O-CNN
Mean IoU	<b>86.2</b>	86.0	85.9
Preprocess	0 ms	0 ms	86.9 ms
Inference	50.7 ms	241.8 ms	5.3 ms
Postprocess	0 ms	0 ms	842.4 ms

Table 1: Results on ShapeNet Part.

	mIoU
PVCNN (on S3DIS)	<b>54.33</b>
w/o scale normalization	52.62
	mIoU
PVCNN (on ShapeNet)	<b>86.2</b>
devoxelization w/o TI	85.7
feature fusion by concat	86.1
voxel resolution ( $0.75\times$ )	85.7
voxel resolution ( $1.25\times$ )	86.1
1 voxel-conv per block	85.6
3 voxel-convs per block	86.1

Table 2: Ablation studies.

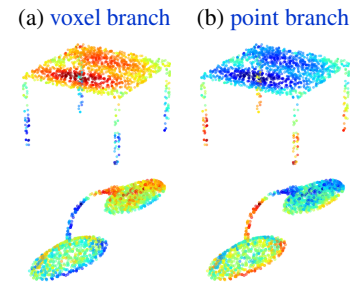


Figure 1: Feature visualization.

