# MASC: Multi-scale Affinity with Sparse Convolution
# for 3D Instance Segmentation
## Technical Report

Chen Liu
Washington University in St. Louis
chenliu@wustl.edu

Yasutaka Furukawa
Simon Fraser University
furukawa@sfu.ca

## Abstract

*We propose a new approach for 3D instance segmentation based on sparse convolution and point affinity prediction, which indicates the likelihood of two points belonging to the same instance. The proposed network, built upon submanifold sparse convolution [3], processes a voxelized point cloud and predicts semantic scores for each occupied voxel as well as the affinity between neighboring voxels at different scales. A simple yet effective clustering algorithm segments points into instances based on the predicted affinity and the mesh topology. The semantic for each instance is determined by the semantic prediction. Experiments show that our method outperforms the state-of-the-art instance segmentation methods by a large margin on the widely used ScanNet benchmark [2]. We share our code publicly at https://github.com/art-programmer/MASC.*

## 1. Introduction

The vision community has witnessed tremendous progress in 3D data capturing and processing techniques in recent years. Consumer-grade depth sensors enable researchers to collect large-scale datasets of 3D scenes [2, 1, 7]. The emergence of such datasets empowers learning-based methods to tackle a variety of 3D tasks, such as object recognition, part segmentation and semantic segmentation. Among them, 3D instance segmentation is an important yet challenging task as it requires accurate segmentation of 3D points based on instances without a fixed label set. In this paper, we propose a simple yet effective method to learn the affinity between points based on which points are clustered into instances.

The irregularity of 3D data poses a challenge for 3D learning techniques. Volumetric CNNs [21, 15, 12] are first explored as they are straightforward extensions of their successful 2D counterparts. Various techniques are proposed to reduce the cost of expensive 3D convolutions. OctNet [17] and O-CNN [19] utilize the sparsity of 3D data via the octree representations to save computation. Other methods [9, 14] process 3D point cloud directly without voxelization. While showing promising results, these methods lack the ability of modeling local structures of the input point cloud. Later methods [16, 8, 10] model local dependencies with various tricks, but the number of points processed by the network is still quite limited (e.g., $4,096$). A point cloud of an indoor space usually contains much more number of points, which means the network can only process a slice of the input point cloud at each time, which disables global reasoning of the space. Recently, Graham *et al.* [3] propose an super-efficient volumetric CNN based on submanifold sparse convolution [4] to process the entire point cloud of an indoor scene, which achieves promising results on the semantic segmentation task. In this paper, we adopt sparse convolution and propose a clustering algorithm based on learned multi-scale affinities to tackle the 3D instance segmentation problem.

Instance segmentation is more challenging than semantic segmentation as it requires the additional reasoning of objects. Instance segmentation methods can be categorized into two groups, proposal-based approaches and proposal-free approaches. Proposal-based approaches builds system upon object detection and append segmentation modules after bounding box proposals. Inspired by the recent success of Mask R-CNN [5] on 2D instance segmentation, 3D-SIS [6] develops a proposal-based system which achieves the state-of-the-art performance on 3D instance segmentation evaluated using the ScanNet benchmark [2]. GSPN [22] presents a generative model for generating proposals. FrustumNet [13] un-project 2D proposals to 3D space for the segmentation network to process. On the other hand, proposal-free methods cluster points into instances based on the similarity metrics. SGPN [20] trains a network to predict semantic labels, a similarity matrix between all pairs of points, and point-wise confidence for being a seed point, from which instance segmentation results are generated. However, the similarity between most pairs of points
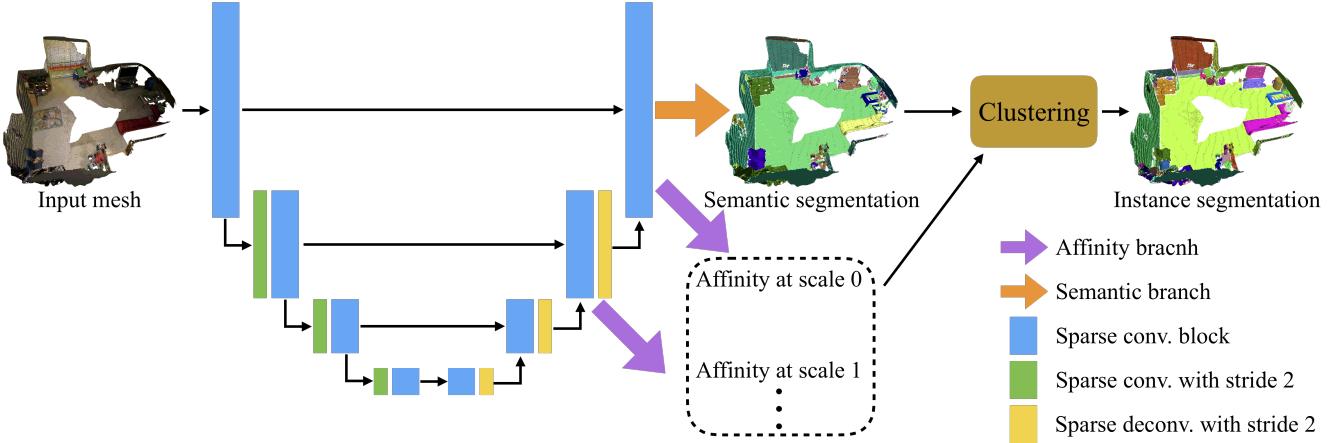
Figure 1. We use a U-Net architecture with submanifold sparse convolutions [3] to process the entire point cloud of an indoor scene and predict semantic scores for each point as well as the affinity between neighboring voxels at different scales. A simple yet effective clustering algorithm groups points into instances based on the predicted affinity and the mesh topology.

is not informative and introduces unnecessary challenge, limiting the network to process only $4,096$ points at each time. We address this issue by predicting only similarity between neighboring points at multiple scales, as did in [11] for 2D instance segmentation, and develops a clustering algorithm to group points based on the learned local similarity and the input mesh topology.

## 2. Methods

As shown in Fig. 1, we use the same U-Net architecture [18] with submanifold sparse convolution used in [3] for semantic segmentation. The sparse U-Net first voxelized the entire point cloud from each ScanNet scene, with each point represented by its coordinate, color, and local surface normal. Following [3], we set the voxel size as $2cm \times 2cm \times 2cm$ and use $4,096 \times 4,096 \times 4,096$ voxel grids so that the entire scene can be voxelized sparsely and each voxel usually contains at most 1 point. After voxelization, submanifold sparse convolution layers with skip connections generate feature maps of different spatial resolutions. We append a semantic branch with one fully connected layer to predict point-wise semantics and add multiple affinity branches at different scales to predict the similarity scores between an active voxel and its 6 neighbors. We denote the finest scale, which has $4,096 \times 4,096 \times 4,096$ voxels, as scale 0 and scale $s$ has resolution $\frac{4,096}{s^2} \times \frac{4,096}{s^2} \times \frac{4,096}{s^2}$.

We treat the input mesh as the initial graph $(V, E)$, and propose a clustering algorithm to group points into instances based on the multi-scale affinity field and the semantic prediction. Note that the clustering is conducted on the mesh graph while the network predicts the affinity between neighboring voxels with fixed neighborhood. So we first compute the affinity between two notes by taking the average affinity between neighboring voxel pairs con-

---

**Algorithm 1** The clustering algorithm based on multi-scale affinity

---

**Require:** Voxel affinity $A^s(p, q)$ at each scale, input mesh $(V, E)$
  **return** Clustering result for $V$
  **repeat**
    Initialize $A(V_i, V_j)$
    **if** $(V_i, V_j) \in E$ **then**
      $A(V_i, V_j) \leftarrow \text{avg}_s(\text{avg}_{p \in V_i, q \in V_j}(A^s(p, q)))$
    **else**
      $A(V_i, V_j) \leftarrow 0$
    **end if**
    Map $V_i$ to a neighbor $M(V_i)$ with high similarity
    **if** $\max_{V_j}(A(V_i, V_j)) > 0.5$ **then**
      $M(V_i) \leftarrow \text{argmax}_{V_j}(A(V_i, V_j))$
    **else**
      $M(V_i) \leftarrow V_i$
    **end if**
    Assign $V_i$ to cluster $C_k$ if $M(V_i) \neq V_i, M(V_i) \in C_k$
    Update $E \leftarrow \{(C_k, C_l) | \exists V_i, V_j : V_i \in C_k, V_j \in C_l, (V_i, V_j) \in E\}$
    Update $V \leftarrow C$
  **until** The update does not change anything

---

necting two nodes. At the beginning, each node contains only one point in the original point cloud which occupies one voxel and the affinity between two nodes is simply the affinity between corresponding voxels. As the node grows, a node could occupy multiple voxels in the finest scale, and if no less than $4^s$ points of the node fall inside the voxel at a higher scale $s$, we say that the node occupies this voxel. With node affinities, we map each node $V_i$ to its most similar neighbor $M(V_i)$ if the similarity between them is higher

than 0.5, and to itself otherwise. With the mapping, we can cluster nodes into groups such that every node in a group is mapped to another node in the same group and none of the node in a group is mapped to other groups. We treat each cluster as a new node and update edges correspondingly. We repeat the process until no change happens (i.e., every node is mapped to itself). The above procedure is summarized in Alg. 1.

Compared against the clustering algorithm in [11] for 2D instance segmentation, our clustering algorithm is more aggressive as it merges nodes in parallel, and has the potential of being implemented using GPU operations. After the clustering process, each instance takes the semantic label with the maximum votes from its points.

# 3. Results

## 3.1. Implementation details

We implement the proposed method using PyTorch based on the code of [3][1]. Besides voxelizing and augmenting ScanNet scenes in the same way of [3], we add point normals to the input and use instance labels to generate affinity supervision between neighboring voxels. Two neighboring voxels at scale $s$ have a similarity score 1 if 1) each of them contains at least $4^s$ points and 2) their instance distributions among points are the same. To better predict local affinity at regions where the input mesh is sparse, we further augment the input mesh by randomly sampling 5 points inside each triangle which spans at least 2 voxel at either dimension. An edge is added between two sampled points if they are voxelized into neighboring voxels at the finest scale, and a sampled point is also connected with its closest vertex of the original triangle. This random densification also enriches the training data. We use 2 scales for the affinity prediction. We train the network on a Titan 1080ti GPU for 20 epochs.

## 3.2. Quantitative evaluation

We evaluate the instance segmentation performance on the ScanNet benchmark [2]. To measure the confidence of each instance for the evaluation purpose, we collect all predicted instances in the training set and train a smaller network to predict if an instance is valid for the predicted label. The network encodes the point cloud of the instance using sparse convolutions and encodes the predicted label with a fully connected layer. We concatenate two features and add two fully connected layers to predict the final confidence score. In addition, we find that the clustering algorithm sometimes struggles to segment objects which are co-planar with other larger objects. So we add additional instances for semantic classes which are often planar, namely

---

[1]https://github.com/facebookresearch/SparseConvNet

pictures, curtains, shower curtains, sinks, and bathtubs (usually contain only the bottom surfaces), by finding connected components of such classes based on the predicted semantics. Our method out-performs the state-of-the-art methods by a large margin. Table 1 shows the comparison against published results (with citable publications).

## 3.3. Qualitative evaluation

We show qualitative results for ScanNet scenes in the validation set in Fig. 2.
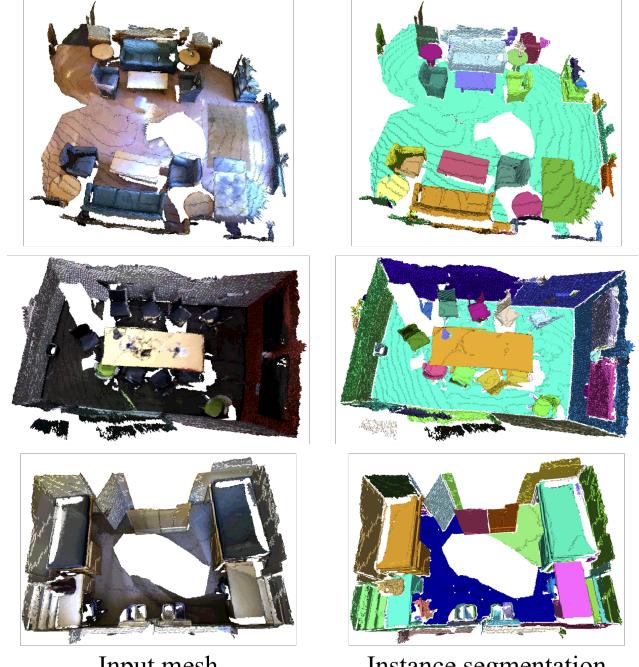


Input mesh          Instance segmentation

Figure 2. Qualitative results for ScanNet scenes. In each row, we show the input mesh on the left and our instance segmentation result on the right.

# 4. Discussion

Though this simple method already achieves promising results, several improvements are required to further boost its performance. First, the clustering algorithm is currently implemented sequentially and thus slow. The algorithm is parallel in theory and it is possible to implement it on GPU and extend it to enable back-propagation for end-to-end training. Besides the speed issue, the effect of multi-scale is under-explored. In practice, we find using 2 scales is fast to train and achieves good performance but it is unclear the role played by each scale. With more exploration on the multi-scale affinity, it is possible to design a better clustering algorithm which uses affinity of more scales to achieve better performance. Finally, the current method sometimes fails to distinguish co-planar objects.

Table 1. Instance segmentation evaluation on the ScanNet benchmark (AP with IOU threshold 0.5)

| Method | avg | bath. | bed | book. | cabi. | chair | coun. | curt. | desk | door | other | pict. | refr. | show. | sink | sofa | table | toil. | wind. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3D-SIS [6] | 0.382 | **1.000** | 0.432 | 0.245 | 0.190 | 0.577 | 0.013 | 0.263 | 0.033 | 0.320 | 0.240 | 0.075 | 0.422 | **0.857** | 0.117 | **0.699** | 0.271 | 0.883 | 0.235 |
| GSPN [22] | 0.306 | 0.500 | 0.405 | 0.311 | 0.348 | 0.589 | **0.054** | 0.068 | 0.126 | 0.283 | 0.290 | 0.028 | 0.219 | 0.214 | 0.331 | 0.396 | 0.275 | 0.821 | 0.245 |
| SGPN [20] | 0.143 | 0.208 | 0.390 | 0.169 | 0.065 | 0.275 | 0.029 | 0.069 | 0.000 | 0.087 | 0.043 | 0.014 | 0.027 | 0.000 | 0.112 | 0.351 | 0.168 | 0.438 | 0.138 |
| Ours | **0.447** | 0.528 | **0.555** | **0.381** | **0.382** | **0.633** | 0.002 | **0.509** | **0.260** | **0.361** | **0.432** | **0.327** | **0.451** | 0.571 | **0.367** | 0.639 | **0.386** | **0.980** | **0.276** |

# 5. Acknowledgement

# References

[1] I. Armeni, S. Sax, A. R. Zamir, and S. Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *arXiv preprint arXiv:1702.01105*, 2017.

[2] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, volume 1, page 1, 2017.

[3] B. Graham, M. Engelcke, and L. van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9224–9232, 2018.

[4] B. Graham and L. van der Maaten. Submanifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307*, 2017.

[5] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017.

[6] J. Hou, A. Dai, and M. Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. *arXiv preprint arXiv:1812.07003*, 2018.

[7] B.-S. Hua, Q.-H. Pham, D. T. Nguyen, M.-K. Tran, L.-F. Yu, and S.-K. Yeung. Scenenn: A scene meshes dataset with annotations. In *3D Vision (3DV), 2016 Fourth International Conference on*, pages 92–101. IEEE, 2016.

[8] Q. Huang, W. Wang, and U. Neumann. Recurrent slice networks for 3d segmentation of point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2626–2635, 2018.

[9] R. Klokov and V. Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 863–872. IEEE, 2017.

[10] C. Liu, J. Wu, and Y. Furukawa. Floornet: A unified framework for floorplan reconstruction from 3d scans. *arXiv preprint arXiv:1804.00090*, 2018.

[11] Y. Liu, S. Yang, B. Li, W. Zhou, J. Xu, H. Li, and Y. Lu. Affinity derivation and graph merge for instance segmenta-tion. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 686–703, 2018.

[12] D. Maturana and S. Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *Intelligent Robots and Systems (IROS), 2015 IEEE/RSJ International Conference on*, pages 922–928. IEEE, 2015.

[13] C. R. Qi, W. Liu, C. Wu, H. Su, and L. J. Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 918–927, 2018.

[14] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 1(2):4, 2017.

[15] C. R. Qi, H. Su, M. Nießner, A. Dai, M. Yan, and L. J. Guibas. Volumetric and multi-view cnns for object classification on 3d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656, 2016.

[16] C. R. Qi, L. Yi, H. Su, and L. J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems*, pages 5099–5108, 2017.

[17] G. Riegler, A. Osman Ulusoy, and A. Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3577–3586, 2017.

[18] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.

[19] P.-S. Wang, Y. Liu, Y.-X. Guo, C.-Y. Sun, and X. Tong. O-cnn: Octree-based convolutional neural networks for 3d shape analysis. *ACM Transactions on Graphics (TOG)*, 36(4):72, 2017.

[20] W. Wang, R. Yu, Q. Huang, and U. Neumann. Sgpn: Similarity group proposal network for 3d point cloud instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2569–2578, 2018.

[21] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015.

[22] L. Yi, W. Zhao, H. Wang, M. Sung, and L. Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. *arXiv preprint arXiv:1812.03320*, 2018.