

GRJointNET: Synergistic Completion and Part Segmentation on 3D Incomplete Point Clouds

Yiğit Gürses
Technical University of Munich

Melisa Taşpınar
Bilkent University

Mahmut Yurt
Stanford University

Sedat Özer
Özyeğin University

Abstract—Segmentation of three-dimensional (3D) point clouds is an important task for autonomous systems. However, success of segmentation algorithms depends greatly on the quality of the underlying point clouds (resolution, completeness etc.). In particular, incomplete point clouds might reduce a downstream model’s performance. GRNet is proposed as a novel and recent deep learning solution to complete point clouds, but it is not capable of part segmentation. On the other hand, our proposed solution, GRJointNet, is an architecture that can perform joint completion and segmentation on point clouds as a successor of GRNet. Features extracted for the two tasks are also utilized by each other to increase the overall performance. We evaluated our proposed network on the ShapeNet-Part dataset and compared its performance to GRNet. Our results demonstrate GRJointNet can outperform GRNet on point completion. It should also be noted that GRNet is not capable of segmentation while GRJointNet is. This study¹, therefore, holds a promise to enhance practicality and utility of point clouds in 3D vision for autonomous systems.

Keywords—Point Clouds, Completion, Segmentation.

I. INTRODUCTION

With the new developments in image acquisition technologies and the widespread use of 3D sensors, the demand for 3D object processing algorithms has also increased. Various algorithms have been developed recently for this purpose as in [1], [2], [3], [4], as well as [5], which forms the basis of our method. A particular relevant application is the processing of 3D point clouds which are often obtained with sensors such as Lidar [6]. 3D objects are commonly represented by 3D point clouds, as point clouds can effectively represent and describe the same scene with significantly lower data size [7] compared to voxel based methods. However, 3D point clouds obtained with sensors tend to be incomplete due to various factors such as light reflection, occlusion, low sensor resolution and limited viewing angles [5], [7], [8]. Hence, the performance of algorithms that use the data as it is suffer [9]. For this reason, a pre-processing step that implements some form of completion and resolution enhancement is often included [10]. GRNet [5] is one of the recently proposed deep learning- based algorithms for this 3D point cloud completion.

Part-segmentation is another type of vision task used in various domains, such as [11], [12] and [13], where each point is assigned one of the predefined labels to segment the whole object into smaller meaningful parts. However, segmenting an incomplete object where some parts may be wholly missing can be unproductive. In this context, completion algorithms can be used as an intermediary step before segmentation to

obtain better results [10]. However, such multi-step processes typically require more resources and can not be parallelised, leading to longer run-times. Therefore, a preferable alternative is to perform point cloud completion and segmentation jointly. In this study, we present a new architecture that aims to simultaneously complete and segment 3D incomplete point clouds, and we call this architecture GRJointNET.

GRJointNET makes use of 3D convolutional layers, three differentiable gridding layers (gridding, gridding reverse, and cubic feature sampling) from [5], a novel segmentation reverse gridding layer and a novel synergistic feature sampling method (see Figure 1). In this method, the incomplete regions in the input point clouds are completed and the points in the created point clouds are and segmented simultaneously.

Our main contributions can be summarized as follows:

- While GRNet cannot perform segmentation together with completion, GRJointNet can perform both segmentation and completion synergistically.
- Unlike the GRNet architecture, our GRJointNet architecture uses segmentation estimates while performing incomplete point completion in the last layer.
- Comparative experimental results on the Shape-Net Part dataset are presented.

II. RELATED WORKS

Several recent studies have presented various deep neural network models to segment and complete 3D objects [14], [15], [16], [17]. One of the studies that pioneered point cloud-based research in this field is PointNet [16]. Although this type of 3D point space-based models have demonstrated some success in the segmentation task, their performance depends on the completeness of the points in the point cloud [7]. However, as mentioned previously, 3D point clouds tend to be incomplete for many reasons [5], [7], [8]. In other words, when working on 3D point clouds, during applications such as segmentation, completing the incomplete point clouds first is considered a separate task.

Many recent and independent studies have successfully demonstrated that incomplete point clouds can be completed using deep neural architectures [5], [7], [18], [19]. Some of these studies perform the completion process using multi-layer perceptrons (MLP) on raw point clouds [7]. However, such MLP-based methods have difficulty in exploiting spatial correlations between points due to the context-unaware architecture of MLPs. For this reason, newer studies have aimed to utilize 3D CNN’s (convolutional neural networks) by voxelizing the point clouds. Even so, in such studies, performance decreases

¹This paper is the authors’ enhanced version. This work was originally published at IEEE SIU 2021 and available on IEEE Xplore with DOI:10.1109/SIU53274.2021.9477918.

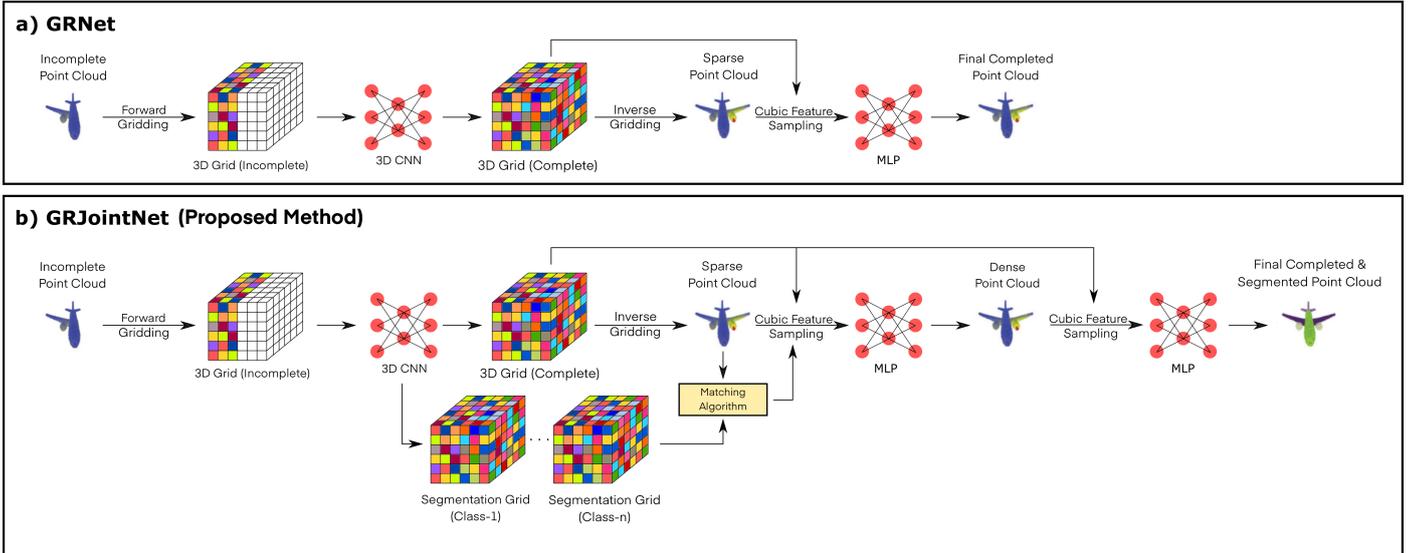


Figure 1: a) The architecture of the base method (GRNet [5]) and b) the architecture of our proposed method are shown. GRJointNet takes an incomplete point cloud as input and processes this cloud data through two different branches (completion and segmentation) to output a completed and segmented point cloud.

may be observed due to loss of geometric information during the voxelization process [20]–[22]. A recent approach, GRNet [5], proposes a model that represents point clouds with 3D grids with the aim of preserving geometric and structural information. Although GRNet is relatively successful at its purpose, it does not have segmentation capabilities.

In this study, we enhance the GRNet structure and present an end-to-end architecture that performs both completion and segmentation simultaneously. We call our architecture GRJointNET. GRJointNET, using GRNet as its base structure, is designed to improve the capabilities of GRNet by incorporating point cloud completion into its framework.

III. THE PROPOSED ARCHITECTURE

The architecture of our proposed method (GRJointNet) is given in Figure 1. In the GRJointNet architecture, there are five fundamental components including (i) gridding, (ii) gridding reverse, (iii) cubic feature sampling, (iv) the 3D convolutional neural network, (v) the multilayer perceptron, (vi) the mapping algorithm and (vii) the loss functions.

Below, we explain each of those components.

1) *Gridding*: It is not defined how to apply 2D and 3D convolutions directly on irregular point clouds, which is why placing the data on a 3D grid structure is a preferred method. Such methods are referred as voxelization. After voxelization, we can apply 2D and 3D convolution operations directly. However, since this process is not reversible, voxelization methods inherently lead to loss of geometric or semantic information. Therefore, in this study, we include a differentiable gridding layer to transform irregular 3D point clouds into regular 3D grids. The targeted 3D grid consists of N^3 individual vertices (where N denotes the number of vertices on one dimension of the grid), covering the entire point cloud given as input and taking the shape of a regular cube. Each cell

in this grid contains 8 different vertices, each with a weight value. The total number of vertices is N^3 with

$$V = \{v_i\}_{i=1}^{N^3}, W = \{w_i\}_{i=1}^{N^3}, v_i = (x_i, y_i, z_i). \quad (1)$$

Here, W holds the cell values whereas the set V holds the vertex coordinates of the corresponding cells. v_i defines the 3D point at the i^{th} index. If a point from the point cloud object lays within a cell with 8 vertices, the weights of these vertices for that point is determined as follows:

$$w_i^p = (1 - |x_i^v - x|)(1 - |y_i^v - y|)(1 - |z_i^v - z|) \quad (2)$$

Here, x represents the projection of a sample coming from the point cloud onto the x-axis, y represents its projection onto the y-axis and z onto the z-axis. x_i^v , y_i^v and z_i^v define a vertex neighbouring the point in question. The final weight w_i of the vertex is then calculated as follows: $w_i = \sum_{p \in N(v_i)} \frac{w_i^p}{|N(v_i)|}$,

where $N(v_i)$ is the set of points neighbouring the vertex v_i . The condition that a point p neighbors v_i can be written as $|x_i^v - x| < 1, |y_i^v - y| < 1, |z_i^v - z| < 1$.

2) *Gridding Reverse*: Gridding reverse is the operation that creates the sparse point cloud from the given 3D grid. The points p_i^s are calculated as follows:

$$p_i^s = \left(\sum_{j \in N(v_i)} w_j v_j \right) / \left(\sum_{j \in N(v_i)} w_j \right), \quad (3)$$

Here, $N(v_i)$ denotes the set of vertices neighboring p_i^s , w_j denotes the weight of the j^{th} vertex in $N(v_i)$, and v_j denotes the spatial position of that vertex.

3) *Cubic Feature Sampling*: Classical MLP-based methods [16] working on 3D point clouds suffer from global and local information loss between neighboring points because they do not take into account local spatial features. To solve this problem, we use the cubic feature sampling technique in our proposed method. This method collects relevant features

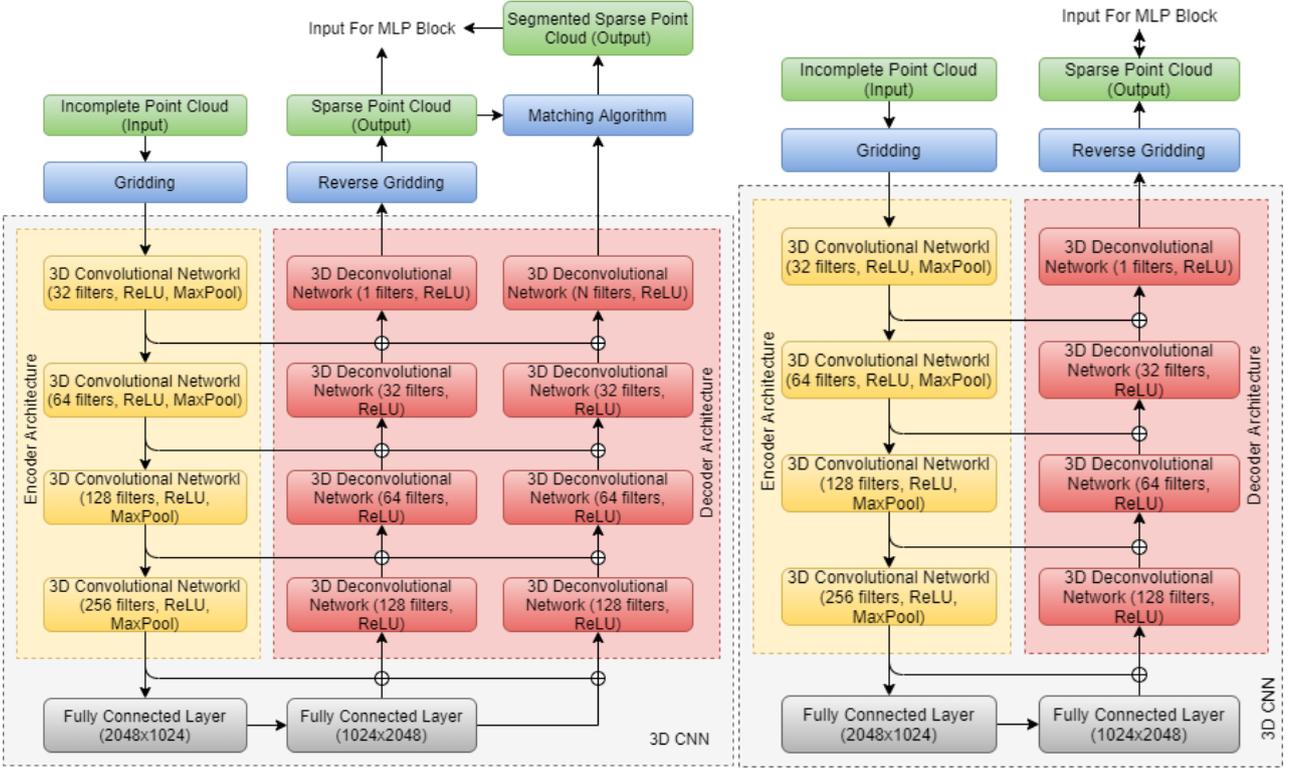


Figure 2: Comparison of the 3D CNN structures of GRJointNet and GRNet. The figure on the left shows the 3D CNN structure used in GRJointNet, and the figure on the right shows the 3D CNN structure used in GRNet.

from the grid for each point in the sparse point cloud. In short, the features of the eight neighboring vertices surrounding the point p_i are combined and the input of the MLP (o_p^i) relative to that point is created as follows: $o_p^i = [p_i, f_1^i, f_2^i, \dots, f_8^i]$. Here, o_p^i denotes the input of the MLP due to the point p_i , whereas f_j^i denotes the feature map of the vertices surrounding p_i from the 3D CNN. Note that the cubic feature sampling takes feature maps from the first three transposed convolutional layers in the 3D CNN, and it randomly samples 8 features from each channel per each point.

4) *3D Convolutional Neural Network*: Both GRNet and GRJointNet each contain a 3D CNN structure. The difference between these two 3D CNN structures can be seen comparatively on Figure 2. The 3D CNN in the proposed approach contains an encoder-decoder structure. The encoder consists of four 3D convolutional layers, each of which includes a padding of 2, batch normalization, max pooling layers of kernel size 4, and a leaky ReLU activation. It is followed by fully connected layers of dimensions 1024 and 2048. Meanwhile, the decoder contains four transposed convolutional layers, each of which includes a padding of 2, stride of 1, a batch normalization, and a leaky ReLU activation. The general formulation of the 3D CNN is defined as follows: $W' = 3DCNN(W)$; where W is the output of the incomplete point cloud from the gridding process, and W' is its completed version. Thus, the 3D CNN recovers the missing points in the given incomplete point cloud.

5) *Multilayer Perceptron (MLP)*: The MLP architecture in the proposed method aims to recover fine details from the

sparse point cloud by using the deviation between the final completed/segmented point cloud and the sparse point cloud. The MLP architecture encompasses four fully connected (FC) layers with sizes 12, 1000, 2000, and 3584, respectively.

6) *Mapping Algorithm*: The performance of GRJointNet depends on the efficient use of the deconvolutional layers that form the segmentation grid to learn well. For this purpose, we segmented the sparse point cloud and used this segmentation in back-propagation with cross entropy loss. The mapping algorithm works as follows:

$$c_x^p = \lfloor N(p_x + 1) \rfloor, c_y^p = \lfloor N(p_y + 1) \rfloor, c_z^p = \lfloor N(p_z + 1) \rfloor, \\ \text{and } b^p = \arg \max_n BI_n[c_x^p, c_y^p, c_z^p]. \quad (4)$$

Here c_x^p, c_y^p and c_z^p indicate the indices of the cell that point p will fall into in a segmentation grid of size N^3 . BI_n denotes the n^{th} of the resulting n segmentation grids and contains the spatial probabilities of the segmentation category numbered n . b^p indicates the segmentation category assigned to point p at the end of the mapping algorithm.

7) *Loss Functions*: The *Chamfer distance* between the actual ground truth and the completed/segmented objects is defined as:

$$L_{CD} = \frac{1}{n_G} \sum_{g \in G} \min_{m \in M} \|g - m\|_2^2 + \frac{1}{n_M} \sum_{m \in M} \min_{g \in G} \|g - m\|_2^2 \quad (5)$$

For each point in G , the closest point in M is calculated based on the distance L_2 . This L_2 distance is included in the loss.

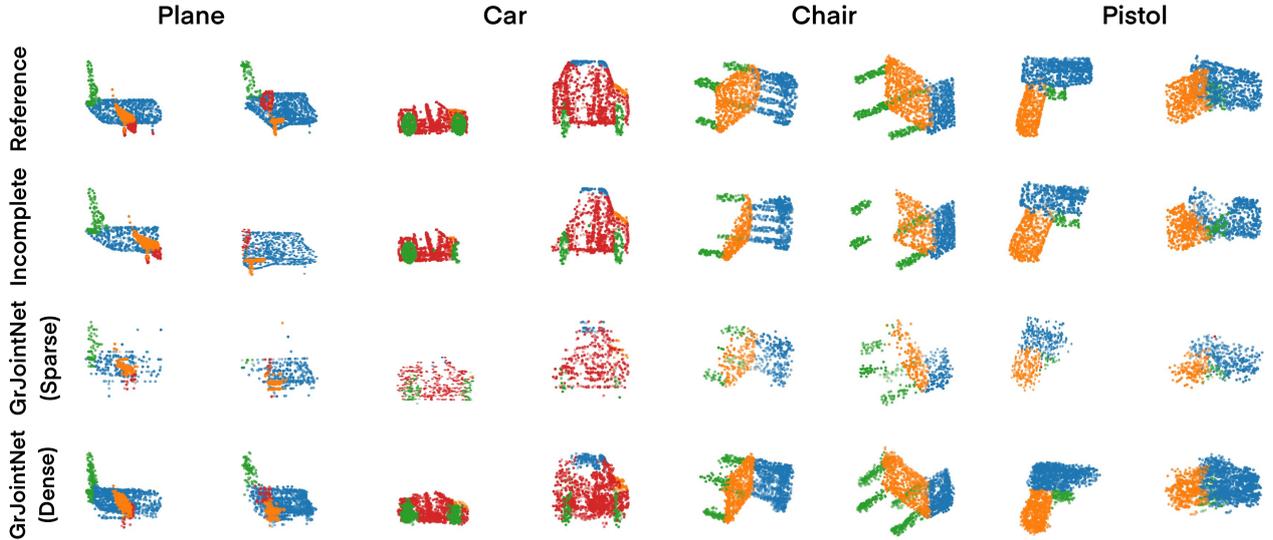


Figure 3: The results of incomplete point cloud completion performed by the proposed GRJointNet model are shown on incomplete, sparse and dense point clouds. Two samples were used for each of the following categories: plane, car, chair and pistol.

The same process is repeated for each point in M . For the segmentation loss, *cross entropy loss* was used.

$$L_{CE} = - \sum_i t_i \log \left(\frac{e^{s_i}}{\sum_j e^{s_j}} \right) \quad (6)$$

Given that complete point clouds do not have ground truths for segmentation when they are first created; the ground truths are calculated using the original complete point cloud which has segmentation labels available. Each point in the generated clouds is assigned to the segmentation label of the point closest to it in the complete cloud. Afterwards, the segmentation predictions on both sparse and dense point clouds are compared to the ground truths we generated using cross entropy. Using only the Chamfer distance as a loss function to train GRNet is insufficient to check whether the predicted points match the geometry of the object. For this reason, networks that use only Chamfer distance tend to give an average shape that minimizes the distance of input and output points. This in turn causes a loss of information regarding the details of the object in question. Since point clouds are unsorted, it becomes difficult to apply L_1 / L_2 loss function or cross entropy directly on them. However, the gridding method introduced by GRNet [5] overcomes this problem by converting unsorted 3D point clouds into 3D grids. Therefore, GRNet introduces a novel loss function called *Grid Loss Function*. This loss function is defined as the distance L_1 between two sets of values of 3D grids. In other words:

$$L_{Gridding}(W^{pred}, W^{gt}) = \frac{1}{N_G^3} \sum ||W^{pred} - W^{gt}||. \quad (7)$$

Here, $W^{pred}, W^{gt} \in \mathbb{R}^{N_G^3}$. $G_{pred} = \langle V^{pred}, W^{pred} \rangle$ and $G_{gt} = \langle V^{gt}, W^{gt} \rangle$ are 3D grids obtained by applying gridding to the ground truth (G_{gt}) and the predicted (G_{pred}) point clouds. Additionally, N_G corresponds to the resolution of the 3D grids. The last used loss function (L) on the other hand is defined as follows: $L = L_{CD} + L_{CE} + L_{Gridding}$.

TABLE I: GRNET VS GRJOINTNET RESULTS

	GRNet	GRJointNet
car	6.26 / 2.92	6.18 / 3.00
plane	5.70 / 1.49	3.27 / 1.50
chair	5.52 / 2.92	4.58 / 2.36
pistol	13.07 / 1.83	12.71 / 1.85

IV. EXPERIMENTS

The performances of GRNet and GRJointNet were compared for four selected categories on the ShapeNet-Part dataset [23]. Given that GRNet is an algorithm designed to perform completion, we carried out our experiments separately for both completion and segmentation purposes. All algorithms were trained over 50 epochs. Adam optimization was used on both networks. In the completion experiments, a total of 11705 training samples and 2768 test samples were used from the ShapeNet-Part dataset. The results are shown comparatively in Table I over four randomly selected individual classes including "car", "plane", "chair" and "pistol". In the table, we used the average Chamfer distance as the metric for performance comparison, where the smaller values are the better results and the best results are shown in bold. For each value pair cd_{sparse}/cd_{dense} in Table I, cd_{sparse} and cd_{dense} refer to the Chamfer distances of the sparse and dense completed point clouds to the ground truth, respectively.

In the part-segmentation experiment, since GRNet does not have a segmentation feature, we present only our results in the Figure 3 using two examples from four different categories. In the figure, the first row shows the reference images, the second row shows the inputted incomplete point clouds, whereas the third and fourth rows respectively show the sparse and dense point clouds that are the outputs of the model, all together with the segmentation results.

V. CONCLUSION

In this study, a synergistic deep learning-based method is proposed for the completion and segmentation of incomplete (3D) point clouds. While the proposed method achieves near or better performance than our baseline method (GRNet) in the completion category, it can also successfully segment the completed point cloud to provide further functionality. In real-world autonomous system applications, the collected data is often incomplete and noisy while including data from several different types of sensors. In this context, it can be said that models focusing only on one task will be less efficient and perform worse compared to integrated systems that process all the available data synergistically. To that end, similar to the method proposed in this study, more useful and effective models that can use various features of the gathered data (position, distance, image, etc.) to perform multiple autonomous system-based tasks are being developed [24]–[26]. Integrated systems using such mentioned methods allow use of common inputs at different components and as such, they can lower the required computational resources, while acquiring extra information from other components' internal processes to enhance each others performances.

ACKNOWLEDGMENT

This paper has been produced benefiting from the 2232 International Fellowship for Outstanding Researchers Program of TÜBİTAK (Project No:118C356). However, the entire responsibility of the paper belongs to the owner of the paper. The financial support received from TÜBİTAK does not mean that the content of the publication is approved in a scientific sense by TÜBİTAK.

REFERENCES

- [1] X. Lai, J. Liu, L. Jiang, L. Wang, H. Zhao, S. Liu, X. Qi, and J. Jia, "Stratified transformer for 3d point cloud segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8500–8509, June 2022. 1
- [2] M. Afham, I. Dissanayake, D. Dissanayake, A. Dharmasiri, K. Thilakarathna, and R. Rodrigo, "Crosspoint: Self-supervised cross-modal contrastive learning for 3d point cloud understanding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 9902–9912, June 2022. 1
- [3] X. Yu, L. Tang, Y. Rao, T. Huang, J. Zhou, and J. Lu, "Point-bert: Pre-training 3d point cloud transformers with masked point modeling," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 19313–19322, June 2022. 1
- [4] C. Zhou, Z. Luo, Y. Luo, T. Liu, L. Pan, Z. Cai, H. Zhao, and S. Lu, "Ptr: Relational 3d point cloud object tracking with transformer," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8531–8540, June 2022. 1
- [5] H. Xie, H. Yao, S. Zhou, J. Mao, S. Zhang, and W. Sun, "Grnet: Gridding residual network for dense point cloud completion," 2020. 1, 2, 4
- [6] Y. Zhang, Q. Hu, G. Xu, Y. Ma, J. Wan, and Y. Guo, "Not all points are equal: Learning highly efficient point-based detectors for 3d lidar point clouds," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 18953–18962, June 2022. 1
- [7] Z. Huang, Y. Yu, J. Xu, F. Ni, and X. Le, "Pf-net: Point fractal network for 3d point cloud completion," 2020. 1
- [8] P. Xiang, X. Wen, Y.-S. Liu, Y.-P. Cao, P. Wan, W. Zheng, and Z. Han, "Snowflake point deconvolution for point cloud completion and generation with skip-transformer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 5, pp. 6320–6338, 2023. 1
- [9] Y. Chen, Y. Li, X. Zhang, J. Sun, and J. Jia, "Focal sparse convolutional networks for 3d object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5428–5437, June 2022. 1
- [10] X. Wen, P. Xiang, Z. Han, Y.-P. Cao, P. Wan, W. Zheng, and Y.-S. Liu, "Pmp-net++: Point cloud completion by transformer-enhanced multi-step point moving paths," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 1, pp. 852–867, 2023. 1
- [11] O. Sahin, F. E. Doğanay, S. Ozer, and C. H. Chen, *Segmentation of COVID-19 Infected Lung Area in CT Scans with Deep Algorithms*, ch. Chapter 2.1, pp. 33–48. 2022. 1
- [12] D. Li, G. Shi, J. Li, Y. Chen, S. Zhang, S. Xiang, and S. Jin, "Plantnet: A dual-function point cloud segmentation network for multiple plant species," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 184, pp. 243–263, 2022. 1
- [13] K. Pasupa, P. Kittiworapanya, N. Hongngern, and K. Woraratpanya, "Evaluation of deep learning algorithms for semantic segmentation of car parts," *Complex and Intelligent Systems*, vol. 8, 05 2021. 1
- [14] L. P. Tchapmi, C. B. Choy, I. Armeni, J. Gwak, and S. Savarese, "Segcloud: Semantic segmentation of 3d point clouds," 2017. 1
- [15] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 922–928, 2015. 1
- [16] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," 2017. 1, 2
- [17] A. Abbasi, S. Kalkan, and Y. Sahillioğlu, "Deep 3d semantic scene extrapolation," *The Visual Computer*, vol. 35, 02 2019. 1
- [18] M. Liu, L. Sheng, S. Yang, J. Shao, and S.-M. Hu, "Morphing and sampling network for dense point cloud completion," 2019. 1
- [19] T. Groueix, M. Fisher, V. Kim, B. Russell, and M. Aubry, "Atlasnet: A papier-mâché approach to learning 3d surface generation," *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 02 2018. 1
- [20] A. Dai, C. R. Qi, and M. Nießner, "Shape completion using 3d-encoder-predictor cnns and shape synthesis," 2017. 2
- [21] X. Han, Z. Li, H. Huang, E. Kalogerakis, and Y. Yu, "High-resolution shape completion using deep neural networks for global structure and local geometry inference," 2017. 2
- [22] Z. Wang and F. Lu, "Voxsegnet: Volumetric cnns for semantic part segmentation of 3d shapes," 2018. 2
- [23] L. Yi, L. Guibas, V. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, A. Lu, Q. Huang, and A. Sheffer, "A scalable active framework for region annotation in 3d shape collections," *ACM Transactions on Graphics*, vol. 35, pp. 1–12, 11 2016. 4
- [24] D. Gozen and S. Ozer, "Visual object tracking in drone images with deep reinforcement learning," in *2020 25th International Conference on Pattern Recognition (ICPR)*, pp. 10082–10089, 2021. 5
- [25] B. M. Albaba and S. Ozer, "Synet: An ensemble network for object detection in uav images," 2020. 5
- [26] S. Özer, M. Ege, and M. A. Özkanoglu, "Siamesefuse: A computationally efficient and a not-so-deep network to fuse visible and infrared images," *Pattern Recognition*, vol. 129, p. 108712, 2022. 5