

# HPGNN: Using Hierarchical Graph Neural Networks for Outdoor Point Cloud Processing

Arulmolivarman Thieshanthan<sup>\*†</sup>, Amashi Niwarthana<sup>\*†</sup>, Pamuditha Somarathne<sup>\*†</sup>,  
Tharindu Wickremasinghe<sup>\*†</sup>, Ranga Rodrigo<sup>\*</sup>

<sup>\*</sup>Department of Electronic and Telecommunication Engineering, University of Moratuwa, Sri Lanka

**Abstract**—Inspired by recent improvements in point cloud processing for autonomous navigation, we focus on using hierarchical graph neural networks for processing and feature learning over large-scale outdoor LiDAR point clouds. We observe that existing GNN based methods fail to overcome challenges of scale and irregularity of points in outdoor datasets. Addressing the need to preserve structural details while learning over a larger volume efficiently, we propose Hierarchical Point Graph Neural Network (HPGNN). It learns node features at various levels of graph coarseness to extract information. This enables to learn over a large point cloud while retaining fine details that existing point-level graph networks struggle to achieve. Connections between multiple levels enable a point to learn features in multiple scales, in a few iterations. We design HPGNN as a purely GNN-based approach, so that it offers modular expandability as seen with other point-based and Graph network baselines. To illustrate the improved processing capability, we compare previous point based and GNN models for semantic segmentation with our HPGNN, achieving a significant improvement for GNNs (+36.7 mIoU) on the SemanticKITTI dataset. <sup>\*</sup>

## I. INTRODUCTION

We focus on feature learning in outdoor LiDAR point clouds, an important perception manner commonplace in vision based autonomous navigation. A setup for collecting sequential scans to form point clouds includes a 3D LiDAR scanner mounted on a vehicle [1]. One challenge of such large-scale outdoor LiDAR point clouds is the high volume of points; generally millions of points per frame of observation [2], [3]. Computing features for each point, with redundancies at dense regions of the cloud, poses a significant computational burden. Another challenge is the irregular distribution of points. Due to the spreading of the detecting beam from the LiDAR source, the observed cloud is denser in the immediate neighbourhood, and sparser with radial distance from the source. In processing such unordered points, repeated grouping and sampling of large clouds add to the complexity. To alleviate such challenges, a suitably down-sampled graphical representation to learn point features through a Graph Neural Network (GNN) is an idea explored [4]. At each node, these features are refined to improve the semantic feature representation. During a GNN iteration, the receptive field of the point grows, and its feature is updated from the aggregation of point features in the receptive field [5]. The problem with this approach is its poor scalability for large outdoor point clouds. We incorporate the idea of hierarchical

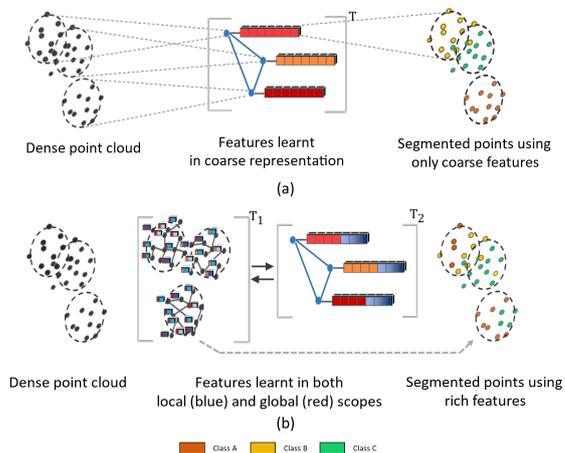


Fig. 1. Acquiring rich features using two levels of coarseness. a) A coarsening approach learns features at one scale, using downsampled points. Features lose fine information. b) A mixture of coarse and fine features and passing between the levels to create rich features. HPGNN follows this approach to segment with higher resolution.

learning and feature pyramid schemes [6] to avoid the need of many iterations to expand the receptive field.

To make feature learning over a large graph computationally efficient, voxelization [4], [7] is a common approach. Such methods vary the sampling scheme based on the structure of the observed point cloud. As shown in Figure 1 (a) they form a coarse level graph representation, learn features from the GNNs, and interpolate them onto the original cloud.

Graph coarsening strategies for reducing computation must strike a balance by sacrificing the accuracy of segmentation [8]. Achieving an optimum coarseness for downsampling is a hard task, since it depends heavily on the dataset and the scale of the objects that we wish to segment. Due to these reasons, Graph Network based attempts that have been made to process large scale outdoor point clouds at the point-scale have thus far not been on par with other approaches.

To address these challenges, we propose a GNN representation—Hierarchical Point GNN (HPGNN)—that can be extended with modularity for feature encoding, and point-cloud processing tasks. We use scalable voxelizations to form graph connections at varying levels of coarseness (Figure 2). We extend the ideas of downsampling and graph coarsening [7] to address the scaling of graph structured data for large outdoor point clouds.

<sup>\*†</sup> These authors contributed equally to this work.

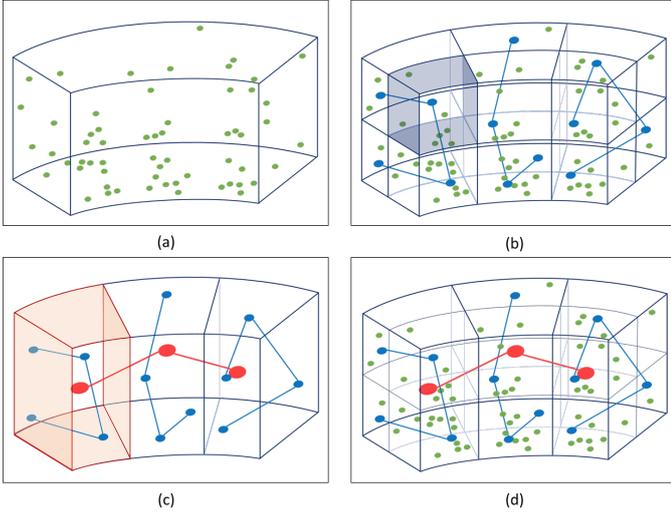


Fig. 2. Graph construction by cylindrical partitioning. The dense point cloud (a) is voxelized at two levels to form a (blue) dense graph (b) and a (red) coarse graph (c). The two resulting GNNs are linked to pass information (d).

HPGNN preserves finer details, while allowing to learn across a larger scale. It leverages existing proximity relations between the points and their features to learn both contextual information locally, as well as spatial structural information in a global scope. As shown in Figure 1 (b), given sufficient learning iterations (T1, T2) at each scale, connections between two hierarchical levels enable a more rich feature representation.

Our main contribution is incorporating a hierarchical learning scheme to improve semantic segmentation of large scale outdoor point-clouds, with computational feasibility in few iterations. We conduct experiments on SemanticKITTI [2] and nuScenes-Lidarseg [9] outdoor datasets to evaluate the proposed architecture. HPGNN outperforms other GNN-based frameworks and other point-based schemes in single scan semantic segmentation.

## II. HPGNN FOR SEMANTIC SEGMENTATION

### A. Voxelization and Keypoint Selection

A common approach to reduce the point cloud density before processing is to use a voxel grid. PointGNN [4], SSCN [10], SegCloud [11] and GridGCN [12], use cubic voxel grids. They do not exploit the irregular distribution within a point cloud. Zhu *et al.*[7] propose a cylindrical voxel partition with asymmetric 3D convolutions as an improvement on regular voxel-based models. These methods voxelize the clouds at a single scale. As the point cloud gets larger and irregular, voxelization schemes to prevent abstracting out fine details are not explored.

We exploit the inherent cylindrical nature of LiDAR data following [7], to parametrise the 3D space using cylindrical coordinates  $(r, \theta, z)$ . We create a scalable voxelization scheme where linear intervals of adjustable magnitude along the radial distance, azimuth angle, and vertical height define

voxel boundaries. This creates larger voxels as the radial distance increases. In a sense, the voxel size defines the “sieve” through which we filter the highly dense LiDAR points. The voxelization is parametrised such that each cylindrical block corresponds to a representative “key point”. A key point is a point with the mean position of the points inside a given voxel. The key points thus formed will create the nodes of the subsequent graph.

Consider two cylindrical grids parametrised by  $(r_L, \theta_L, z_L)$  and  $(r_H, \theta_H, z_H)$ ; each with a different coarseness to sample the point cloud (Figure 2 (b), (c)). They form nodes of two graphs, corresponding to two levels of coarseness (Figure 2 (d)). We maintain connections between the two graphs so that each point in a finer layer may learn from a wider scope through its connection to the coarser layer. This enables to process a wider scope of the point cloud, while preserving structural information that would otherwise be lost through down-sampling. This idea may be extended to create more than two levels in the HPGNN.

### B. Hierarchical Point Graph

Within the HPGNN, a graph  $G = (P, E)$  would be constructed with  $N$  points, where  $P = \{p_1, p_2, \dots, p_N\}$  are the set of keypoints sampled at a particular level of coarseness. Considering two adjacent levels of coarseness, we define the low-level graph  $G_L = (P_L, E_L)$ , and high-level graph  $G_H = (P_H, E_H)$ .  $P_L$  is formed from the smaller voxel size to learn fine details related to a locality.  $P_H$  contains keypoints from larger voxels, to learn more global and structural features from the network. Each point is characterised by  $p_u = (s_u, x_u)$ , with a location vector  $x_u$  and a feature vector  $s_u$ . Hyperparameters  $d_L$  and  $d_H$  are radii defining the edges between points  $p_u, p_v$  that form the set of edges  $E_i$  for each level  $i \in (L, H)$ .

$$E_i = \{(p_u, p_v) \mid \Delta_x < d_i\}, \quad \Delta_x = \|x_u - x_v\|_2. \quad (1)$$

The two graphs  $G_L$  and  $G_H$  are connected based on the voxels that define their points. Each point  $q$  in  $P_L$  will connect with a point  $p$  in  $P_H$ , if the voxel of  $p$  encapsulates  $q$ . After defining two graph levels and the spatial connection between them, we construct a formalism to represent how learning is possible through these connections, with the “neural message passing” approach by Hamilton *et al.*[5], [13]. A node of the network has the aim of progressively refining its feature. In each iteration  $t$ , it uses its current feature vector and relative coordinates to generate a message/edge feature  $e$  through a learnable message generating function  $f(\cdot)$ .

$$e_{uv}^t = f(s_u, s_v, \Delta_x). \quad (2)$$

The destination vector aggregates the messages that are received from its neighbours using an aggregator  $\text{Agg}(\cdot)$  and updates its feature vector  $s$  through a learnable function  $g(\cdot)$ . A separate aggregator converts features between  $G_L$  and  $G_H$  (down/up sampling).

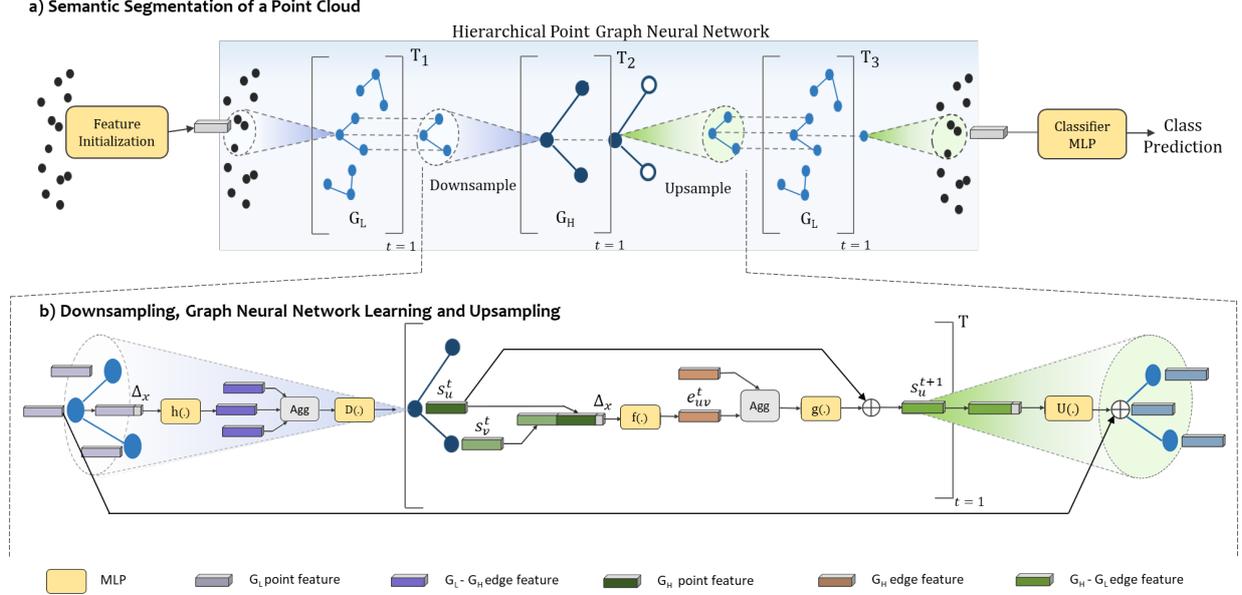


Fig. 3. a) Overview of the HPGNN architecture with feature learning in the low-level graph ( $G_L$ ) and the high-level graph ( $G_H$ ). b) A downsampling layer (blue) that maps local features extracted from points in  $G_L$  to  $G_H$ , and an upsampling layer (green) that maps the learnt features from  $G_H$  back to  $G_L$  for classification. A skip connection adds the input feature for the three layer block (Downsampling,  $G_H$ , and Upsampling) to its output.

### C. Learning in the Hierarchical Graph

Through design decisions, we strive for better GNN learning through an architecture that supports global and local message passing. Furthermore, an appropriate choice of aggregation functions, and strategies to handle neighbourhoods of varying density and similarity are considered. The following subsections refer to Figure 3 and describe the learning process within two adjacent levels, which could be generalized for  $n$  layers.

**Lower Graph:** The GNN of  $G_L$  has higher density by design, and most neighbours of a point correspond to the same object in the dataset. Therefore, most edges of  $G_L$  connect nodes of the same label, giving a high homophily as described by Zhu *et al.*[14]. In  $G_H$ , two neighbour nodes usually correspond to different labels, and have less homophily. We use the Max function as the aggregator  $\text{Agg}(\cdot)$  of features from neighbouring nodes. This has the advantage of not being prone to an over-smoothing which is common in Graph Convolution Networks with mean aggregation [15]. The use of Max function should be complemented by a preprocessing step of removing outlier points, to reduce the possibility of outlier features being aggregated into a node feature. After a feature is aggregated, the conventional approach is to use a multi-layer perceptron (MLP) to learn the function  $g(\cdot)$  and finally add  $s^t$  in each iteration  $t$ .

$$s_u^{t+1} = g(\text{Agg}\{e_{uv}^t | (u, v) \in E_i\}) + s_u^t \quad (3)$$

for each level  $i \in (L, H)$ . This creates a skip connection to ease learning and gradient flow as the networks get deeper. Specifically for GNN, it encourages a point feature not to deviate too far from its previous feature.

After  $T_1$  iterations of message passing,  $G_L$  has refined its

features locally. Then points in  $G_H$  aggregate  $G_L$  features to initialise  $P_H$  node features through a "downsampling" layer. Learning in  $G_H$  continues for  $T_2$  iterations, after which  $P_H$  points return the features to the corresponding  $P_L$  points through an "upsampling" layer.

**Downsample:** Consider a point  $p = (s_p, x_p) \in P_H$  and the corresponding set of points  $\{q_1, q_2, \dots, q_i, \dots, q_k\} \in P_L$  through which edges between  $G_H$  and  $G_L$  are formed. Each  $q_i = (s_{q_i}, x_{q_i})$  has a learnt feature  $s_{q_i}$  from the preceding  $T_1$  iterations. This is concatenated with  $\Delta_x = \|x_p - x_{q_i}\|$ , and a function  $h(\cdot)$  learns the edge feature  $e_{p, q_i}$  between  $p$  and  $q_i$ .

$$e_{p, q_i} = h(s_{q_i}, \Delta_x) \quad (4)$$

The downsampled feature to  $G_H$  is the aggregate of such edge features between the point  $p$  and  $q_i$ ,  $i \in [1, k]$ . For aggregation, we use attention weights for each edge feature following an Attentive Feature Merging (AFM) scheme [16].  $\alpha_i$  is the attention weight for  $e_{p, q_i}$ . We give more representational power for the aggregation stage in the downsampling layer, since a feature representation of a larger volume should reflect the less homophilic nature of  $G_H$ . The aggregated edge feature is then transformed through the MLP for the downsampling function  $D(\cdot)$ .  $s_p^0$  is the initial feature of the point  $p \in P_H$ .

$$s_p^0 = D(\sum(\alpha_i e_{p, q_i})) \quad (5)$$

**Higher Graph:** An almost exact sequence of message passing and aggregating, initialised by features in the method shown by Eq.5 is used. When compared to  $G_L$ , the difference is the number of iterations  $T_2$  that we use for the higher level. If

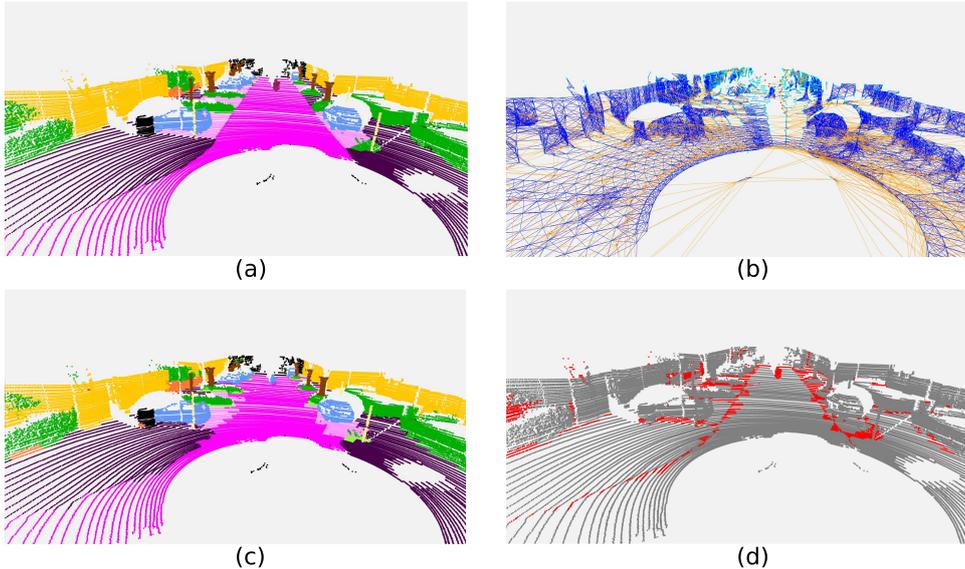


Fig. 4. (a) Ground truth labelled point cloud. (b) Construction of a HPGNN at two levels. The lower level and higher level are blue and orange respectively. (c) Segmentation output of the 2-level HPGNN. (d) Difference between the labels and predictions.

■ unlabelled ■ car ■ motor-cyclist ■ road ■ traffic-sign ■ sidewalk ■ building ■ parking ■ fence ■ pole ■ vegetation ■ terrain ■ trunk

the HPGNN has more than two levels, as the level of  $G_H$  gets higher, two neighbour nodes increasingly correspond to different labels, and have less homophily.

**Upsample and Re-iteration:** When  $G_H$  has completed iterating and learning over a wide neighbourhood, the learnt global message of each point  $p = \{s_p, x_p\} \in P_H$  is passed down as an input to the upsampling function  $U(\cdot)$ . The output feature is concatenated with the  $G_L$  feature  $s_{q_i}$  at the start of the downsample layer, to form a skip connection that skips the 3-layer block of a downsample layer, a GNN, and an upsample layer. Having both higher and lower level information at the starting feature vector, the lower GNN re-iterates  $T_3$  steps to refine its features. We use a similar downsampling to initialise  $G_L$  features from the initialised point features from the original LiDAR cloud. A similar upsampling is used to map refined  $G_L$  features back to the original points. Finally, a classifier MLP predicts classes for each point.

**Loss Function:** In the classifier MLP, we predict a multi-class probability distribution for each point. Since large point cloud data sets have class distributions that are severely imbalanced, we use a weighted cross entropy loss  $L_{wce}$  following [24]. For each class  $c$ , a weight  $w$  is defined inversely proportional to the relative frequency of the class. When  $w_c$  and  $\hat{y}_i$  are the weight and predicted probability of the labelled class, at a point  $i$  in a point cloud with  $N$  points, we have

$$L_{wce} = -\frac{1}{N} \sum_{i=1}^N w_c \ln(\hat{y}_i). \quad (6)$$

To further improve the mean Intersection-over-Union (mIoU) of the model using the Lovasz extension of the Jaccard loss, we follow [25] by using the Lovasz-softmax loss  $L_{ls}$ . Adding an  $L_2$  regularisation term  $L_{reg}$ , the total loss for semantic

segmentation is

$$L_{total} = \alpha L_{wce} + \beta L_{ls} + \gamma L_{reg}. \quad (7)$$

where  $\alpha, \beta, \gamma$  are weights for each loss component.

### III. EXPERIMENTS AND DISCUSSION

To evaluate if the expected improvements in point cloud learning capability is achieved, we choose semantic segmentation as the point cloud processing task and use a 2-level HPGNN with  $(T_1, T_2, T_3) = (1, 2, 1)$  as the hierarchical model. Following our motivation on outdoor point cloud processing, we evaluate on large scale autonomous driving datasets—SemanticKITTI [2] and nuScenes-Lidarseg [33]. Evaluations follow mIoU recommended by [2] as the evaluation metric.

Figure 4 depicts the process of evaluation on large scale point clouds. (a) The point-wise labelled LiDAR scan. (b) An HPGNN, with  $G_L$  and  $G_H$ .  $G_L$  forms most of its edges within object classes, and some edges at boundaries.  $G_H$  forms edges primarily between classes, and facilitates learning over a larger distance in few iterations. The result of segmentation is in (c), with all the classes being identified, but with some edges between classes showing discretization errors.

To compare our results, we identify two broad categories of implementations that have been used for point cloud learning tasks. The first include point-level networks, and the graph-based methods that have used information at that scale. The main motivation in these models is to use an efficient information aggregating mechanism, and not specifically focused only on semantic segmentation. The second category, are the "second generation models" that extend the point network framework to specifically focus on semantic segmentation.

TABLE I

MEAN INTERSECTION OVER UNION (mIoU) SCORES COMPARED WITH SINGLE SCAN SEMANTIC SEGMENTATION AMONG POINT-BASED AND GNN MODELS. OUR PROPOSED MODEL SIGNIFICANTLY OUTPERFORMS EXISTING GRAPH BASED METHOD (SP GRAPH) BY A MARGIN OF 36.7% mIoU, WHILE PERFORMING COMPETITIVELY WITH THE POINT BASED STATE OF THE ART.

Method	mIoU	Car	Bicycle	Motorcycle	Truck	Other-vehicles	Person	Bicyclist	Motorcyclist	Road	Parking	Sidewalk	Other-ground	Building	Fence	Vegetation	Trunk	Terrain	Pole	Traffic-sign
PointNet [17]	14.6	46.3	1.3	0.3	0.1	0.8	0.2	0.2	0.0	61.6	15.8	35.7	1.4	41.4	12.9	31.0	4.6	17.6	2.4	3.7
PointNet++ [18]	20.1	53.7	1.9	0.2	0.9	0.2	0.9	1.0	0.0	72.0	18.7	41.8	5.6	62.3	16.9	46.5	13.8	30.0	6.0	8.9
SPLATNet [19]	22.8	66.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	70.4	0.8	41.5	0.0	68.7	27.8	72.3	35.9	35.8	13.8	0.0
TangentConv [20]	35.9	86.8	1.3	12.7	11.6	10.2	17.1	20.2	0.5	82.9	15.2	61.7	9.0	82.8	44.2	75.5	42.5	55.5	30.2	22.2
P <sup>2</sup> Net [21]	39.8	85.6	20.4	14.4	14.4	11.5	16.9	24.9	5.9	87.8	47.5	67.3	7.3	77.9	43.4	72.5	36.5	60.8	22.8	38.2
RandLA-Net [22]	53.9	<b>94.2</b>	26.0	25.8	<b>40.1</b>	<b>38.9</b>	<b>49.2</b>	48.2	7.2	<b>90.7</b>	<b>60.3</b>	<b>73.7</b>	20.4	<b>86.9</b>	<b>56.3</b>	<b>81.4</b>	61.3	<b>66.8</b>	49.2	47.7
SPGraph [23]	17.4	49.3	0.2	0.2	0.1	0.8	0.3	2.7	0.1	45.0	0.6	34.8	0.6	64.3	20.8	48.9	27.2	24.6	15.9	0.8
<b>HPGNN - Ours</b>	<b>54.1</b>	92.7	<b>33.2</b>	<b>32.2</b>	26.4	29.4	43.3	<b>60.9</b>	<b>19.4</b>	86.3	50.7	67.5	<b>27.6</b>	86.8	55.5	80.7	<b>61.4</b>	64.1	<b>50.6</b>	<b>59.1</b>

TABLE II

mIoU SCORES COMPARED WITH SINGLE SCAN SEMANTIC SEGMENTATION ON SEMANTICKITTI TEST SET. OUR MODEL PERFORMS COMPETITIVELY WITH SECOND GENERATION MODELS EMPLOYING POINT NETWORKS OR GRAPH NETWORKS AS THEIR UNDERLYING ARCHITECTURE.

Method	mIoU	Car	Bicycle	Motorcycle	Truck	Other-vehicles	Person	Bicyclist	Motorcyclist	Road	Parking	Sidewalk	Other-ground	Building	Fence	Vegetation	Trunk	Terrain	Pole	Traffic-sign
MINet [26]	54.3	85.2	38.2	32.1	29.3	23.1	47.6	46.8	24.5	90.5	58.8	72.1	25.9	82.2	49.5	78.8	52.5	65.4	37.7	55.5
PolarNet [27]	54.3	93.8	40.3	30.1	22.9	28.5	43.2	50.2	5.6	90.8	61.7	74.4	21.7	90.0	61.3	84.0	65.5	67.8	51.8	57.5
SqueezeSegv3 [28]	55.9	92.5	38.7	36.5	29.6	33.0	45.6	46.2	20.1	91.7	63.4	74.8	26.4	89.0	59.4	82.0	58.7	65.4	49.6	58.9
Cylinder3D [7]	67.8	97.1	67.6	64.0	<b>59.0</b>	58.6	73.9	67.9	36.0	91.4	65.1	75.5	32.3	91.0	66.5	85.4	71.8	68.5	62.6	65.6
AF2S3Net [29]	69.7	94.5	65.4	<b>86.8</b>	39.2	41.1	<b>80.7</b>	<b>80.4</b>	<b>74.3</b>	91.3	68.8	72.5	<b>53.5</b>	87.9	63.2	70.2	68.5	53.7	61.5	<b>71.0</b>
SPVNAS [30]	67.0	97.2	50.6	50.4	56.6	58.0	67.4	67.1	50.3	90.2	67.6	75.4	21.8	91.6	66.9	86.1	73.4	71.0	64.3	67.3
AMVNet [31]	65.3	96.2	59.9	54.2	48.8	45.7	71.0	65.7	11.0	90.1	<b>71.0</b>	75.8	32.4	92.4	69.1	85.6	71.7	69.6	62.7	67.2
RPVNet [32]	<b>70.3</b>	<b>97.6</b>	<b>68.4</b>	68.7	44.2	<b>61.1</b>	75.9	74.4	73.4	<b>93.4</b>	70.3	<b>80.7</b>	33.3	<b>93.5</b>	<b>72.1</b>	<b>86.5</b>	<b>75.1</b>	<b>71.7</b>	<b>64.8</b>	61.4
<b>HPGNN - Ours</b>	54.1	92.7	33.2	32.2	26.4	29.4	43.3	60.9	19.4	86.3	50.7	67.5	27.6	86.8	55.5	80.7	61.4	64.1	50.6	59.1

TABLE III

mIoU SCORES COMPARED WITH SINGLE SCAN SEMANTIC SEGMENTATION ON NUSCENES-LIDARSEG TEST SET.

Method	mIoU
MINet [26]	56.4
PolarNet [27]	69.4
Cylinder3D [7]	77.2
AMVNet [31]	77.3
SPVNAS [30]	77.4
AF2S3Net [29]	78.3
SPVCNN++ [30]	81.1
<b>HPGNN - Ours</b>	<b>63.8</b>

TABLE IV

ABLATION STUDY ON THE NEED FOR MULTIPLE LEVELS AND NUMBER OF ITERATIONS IN EACH LEVEL.

T1	T2	T3	mIoU
1	0	1	43.3
0	2	0	41.0
1	1	1	<b>51.5</b>
1	2	1	51.6
1	3	1	48.5
2	1	2	53.4
2	2	2	53.1
3	1	3	<b>55.6</b>

### A. Results on SemanticKITTI

Comparisons of mIoU scores on the SemanticKITTI benchmark with point-based and second generation models are shown in Table I and Table II.

**Point-based Models:** Earlier methods (PointNet[17], PointNet++[18]) use MLPs without a predefined underlying structure. More recent models use ideas of GNNs where an underlying structure is generated among points, and some form of information aggregation/pooling processes information from a fixed locality in each iteration (SplatNet[19], TangentConv[20]). These methods perform well for indoor data sets but fail to retain local information, while increasing the receptive field for large outdoor clouds. This is evident from the class-wise mIoU scores in Table I where in pursuit

of widening the receptive field by many iterations, fine information that helps to classify smaller objects are lost. In addition, they suffer from scalability issues for large point clouds as mentioned in Section I. RandLaNet[22] attempts to solve this issue by random sampling points, but it does not solve the issue of information loss.

SPGraph [23] is the highest performing Graph Neural Network model that has previously attempted semantic segmentation. Although previously mentioned point based models use ideas similar to GNNs, they focus on encoding and attentive pooling modules rather than neighborhood feature aggregation mechanisms. From the experiments it can be seen that by using the hierarchical structure for learning, we can retain local information to segment and differentiate between relatively smaller objects (e.g., motor-cyclist, bicyclist, and

pole) as well as larger objects (e.g., car, terrain, and trunk). This shows our design decisions of preserving fine features while graph coarsening has been effective. Therefore, using MLPs in hierarchical levels greatly improves the baseline performance for the Graph-Neural Network implementation of the point-based learning scheme. To this end, we believe we are successful in improving the baseline of point-based GNNs that can be used for large scale point cloud processing.

**Second Generation Semantic Segmentation Models:** These models in Table II focus on the particular task of semantic segmentation and use modules that are specially tuned to improve the performance metrics. They stand upon earlier mentioned point networks by adding range based transformations (SqueezeSeg[34], SalsaNext[35]), adaptive 3D space partitions (AF2S3Net[29], Cylinder3D[7], Polarnet[27]), and multi view fusion, incorporating information such as bird-eye-view and range-images (MINet [26], AMVNet[31], SPVNAS[30], RPVNet[32]). For semantic segmentation, fusion-based approaches currently achieve state-of-the-art results. RPVNet, for example, uses a three-branch fusion network. Instead of a branch with the point graph which uses a simplified PointNet [17], an HPGNN could be used to learn additional features at various degrees of coarseness. Since we observe that HPGNN outperforms PointNet in Table I, we can expect an improvement in second generation models by employing an HPGNN.

It is noteworthy that our model has not used any of the above enhancements specifically tailored for improving semantic segmentation. We have used hierarchical learning to improve the underlying effectiveness of message passing inside a GNN.

### B. Results on nuScenes-Lidarseg

In 2021, nuScenes-Lidarseg was released [9], with point-wise annotations making it suitable for semantic segmentation. Therefore, to our knowledge, point-based models mentioned in Section III-A have not been evaluated in this test bench publicly. The only available comparisons are with the second generation models. They use point-based models in their feature representation step. (e.g., second generation model PolarNet [27] trains a PointNet[17] upon which “ring convolutions” are operated.) Therefore, if our HPGNN, without any additional module, is performing competitively with such second generation models, as seen from Table III, we infer it is an improvement of the underlying point-based model. This validates the generalizability of HPGNN for large scale outdoor LiDAR point clouds.

### C. Ablation Studies

Ablation studies are done using SemanticKITTI with the standard practice [2], [36] of reserving sequence 08 of the dataset for validations.

**Hierarchical Levels:** To validate the need of multiple levels, two models with a single level GNN each are compared with the baseline, which is a 2-level HPGNN. The first model only runs GNN on  $G_L$ :  $(T_1, T_2, T_3) = (1, 0, 1)$  and second model only runs GNN on  $G_H$  graph:  $(T_1, T_2, T_3) = (0, 2, 0)$ . From

Table IV, it can be seen that a 2-level HPGNN on both  $G_L$  and  $G_H$  results in up-to 10.5 mIoU improvement over single level HPGNNs. Clearly, the inability of learning at different scales hinders the performance of the model.

**GNN Iterations:** The impact of giving prominence to each level of a 2-level HPGNN through the number of allowed iterations to pass messages at each scale is examined. We analyse the contribution of increasing the receptive field of a node in  $G_H$  using three models with varying  $T_2$ ;  $(T_1, T_2, T_3) \in \{(1, 1, 1), (1, 2, 1), (1, 3, 1)\}$ , then the same effect in  $G_L$  using models with  $T_2$  constant. Table IV shows that  $T_2 = 2$  has slight improvement over  $T_2 = 1$ , and that  $T_2 = 3$  reverts back to a lower mIoU. This might be due to the receptive field being unnecessarily large compared to the dimensions of the instances in this dataset.

Increasing  $T_1, T_3$  shows improvements of 4.1 mIoU. This suggests that refining and learning features at  $G_L$  significantly improves performance with moderate reception at  $G_H$ .

### D. Limitations of HPGNN

HPGNN is used as a framework for learning features from a point-based representation of a point cloud. For computational feasibility, we use voxelization. As presented by [37], spatially uniform re-sampling has discretization errors, and this is evident from Figure 4 where the errors of segmentation are mostly at the boundaries.

As opposed to graph-resampling methods including Octrees [38], HPGNN requires dense graphs to achieve high resolution, which leads to spatial inefficiency. This is when redundant points of a dense point cloud are stored and processed that do not contribute for evolving new features. This is a limitation in HPGNN in the graph construction step. Although our idea of connecting local points with similar semantics for efficient representation is explored in a different approach by Hypergraphs[39], it is not entirely dependent on point distance.

## IV. CONCLUSION

We introduced a novel framework for learning discriminative features of large-scale outdoor point clouds, through graph neural message passing. The hierarchical point graph neural network (HPGNN) learns features of various levels of coarseness from point clouds. Our experiments show that it is a resource efficient method for expanding the effective receptive field of each point in a graphical structure. HPGNN was designed as a GNN based approach, and it significantly improved the existing baseline of GNN based point cloud segmentation (+36.7 mIoU on SemanticKITTI dataset). Our results suggest an improved GNN based backbone for point cloud processing, which may be extended upon by modular designs or fusion-based approaches for semantic segmentation and other tasks related to point clouds.<sup>†</sup>

Acknowledgement: We thank National Research Council of Sri Lanka for providing computational resources through the grant no. 19-080.

<sup>†</sup>The code will be made available here: <https://git.io/JVhDc>

## REFERENCES

- [1] G. Biao, P. Yancheng, L. Chengkun, and G. Sibó, "Are we hungry for 3d lidar data for semantic segmentation? a survey and experimental study," in *Proceedings of the IEEE Transactions on Intelligent Transportation Systems*, 2020.
- [2] J. Behley, M. Garbade, A. Milioto, J. Quenzel, S. Behnke, C. Stachniss, and J. Gall, "SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences," in *Proc. of the IEEE/CVF International Conf. on Computer Vision (ICCV)*, 2019.
- [3] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 3354–3361.
- [4] W. Shi and R. Rajkumar, "Point-gnn: Graph neural network for 3d object detection in a point cloud," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [5] W. L. Hamilton, "Graph representation learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 14, no. 3, pp. 1–159, 2020.
- [6] G. Zhao, W. Ge, and Y. Yu, "Graphfpn: Graph feature pyramid network for object detection," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 2763–2772.
- [7] X. Zhu, H. Zhou, T. Wang, F. Hong, Y. Ma, W. Li, H. Li, and D. Lin, "Cylindrical and asymmetrical 3d convolution networks for lidar segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 9939–9948.
- [8] Z. Huang, S. Zhang, C. Xi, T. Liu, and M. Zhou, "Scaling up graph neural networks via graph coarsening," *arXiv preprint arXiv:2106.05150*, 2021.
- [9] W. K. Fong, R. Mohan, J. V. Hurtado, L. Zhou, H. Caesar, O. Beijbom, and A. Valada, "Panoptic nusenes: A large-scale benchmark for lidar panoptic segmentation and tracking," *arxiv preprint arXiv:2109.03805*, 2021.
- [10] 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*, 2018, pp. 9224–9232.
- [11] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, and S. Savarese, "Segcloud: Semantic segmentation of 3d point clouds," in *2017 international conference on 3D vision (3DV)*. IEEE, 2017, pp. 537–547.
- [12] Q. Xu, X. Sun, C.-Y. Wu, P. Wang, and U. Neumann, "Grid-gcn for fast and scalable point cloud learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5661–5670.
- [13] W. L. Hamilton, R. Ying, and J. Leskovec, "Inductive representation learning on large graphs," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 1025–1035.
- [14] J. Zhu, Y. Yan, L. Zhao, M. Heimann, L. Akoglu, and D. Koutra, "Beyond homophily in graph neural networks: Current limitations and effective designs," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 7793–7804.
- [15] D. Chen, Y. Lin, W. Li, P. Li, J. Zhou, and X. Sun, "Measuring and relieving the over-smoothing problem for graph neural networks from the topological view," *AAAI 2020 - 34th AAAI Conference on Artificial Intelligence*, pp. 3438–3445, 2020.
- [16] P. Veličković, A. Casanova, P. Liñeira, G. Cucurull, A. Romero, and Y. Bengio, "Graph attention networks," *6th International Conference on Learning Representations, ICLR 2018 - Conference Track Proceedings*, pp. 1–12, 2018.
- [17] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [18] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," *Advances in Neural Information Processing Systems*, vol. 30, pp. 5099–5108, 2017.
- [19] H. Su, V. Jampani, D. Sun, S. Maji, E. Kalogerakis, M.-H. Yang, and J. Kautz, "Splatnet: Sparse lattice networks for point cloud processing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [20] M. Tatarchenko, J. Park, V. Koltun, and Q.-Y. Zhou, "Tangent convolutions for dense prediction in 3d," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [21] S. Li, Y. Liu, and J. Gall, "Projected-point-based segmentation: A new paradigm for lidar point cloud segmentation," *ArXiv*, vol. abs/2008.03928, 2020.
- [22] Q. Hu, B. Yang, L. Xie, S. Rosa, Y. Guo, Z. Wang, N. Trigoni, and A. Markham, "Randla-net: Efficient semantic segmentation of large-scale point clouds," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 108–11 117.
- [23] L. Landrieu and M. Simonovsky, "Large-scale point cloud semantic segmentation with superpoint graphs," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4558–4567.
- [24] I. Alonso, L. Riazuelo, L. Montesano, and A. C. Murillo, "3d-mininet: Learning a 2d representation from point clouds for fast and efficient 3d lidar semantic segmentation," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 5432–5439, 2020.
- [25] M. Berman, A. R. Triki, and M. B. Blaschko, "The lovász-softmax loss: A tractable surrogate for the optimization of the intersection-over-union measure in neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4413–4421.
- [26] S. Li, X. Chen, Y. Liu, D. Dai, C. Stachniss, and J. Gall, "Multi-scale interaction for real-time lidar data segmentation on an embedded platform," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 738–745, 2022.
- [27] Y. Zhang, Z. Zhou, P. David, X. Yue, Z. Xi, B. Gong, and H. Foroosh, "Polarnet: An improved grid representation for online lidar point clouds semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [28] C. Xu, B. Wu, Z. Wang, W. Zhan, P. Vajda, K. Keutzer, and M. Tomizuka, "SqueezeSegv3: Spatially-adaptive convolution for efficient point-cloud segmentation," in *European Conference on Computer Vision*. Springer, 2020, pp. 1–19.
- [29] R. Cheng, R. Razani, E. Taghavi, E. Li, and B. Liu, "2-s3net: Attentive feature fusion with adaptive feature selection for sparse semantic segmentation network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 12 547–12 556.
- [30] H. Tang, Z. Liu, S. Zhao, Y. Lin, J. Lin, H. Wang, and S. Han, "Searching efficient 3d architectures with sparse point-voxel convolution," in *European Conference on Computer Vision*. Springer, 2020, pp. 685–702.
- [31] V. E. Liong, T. N. T. Nguyen, S. Widjaja, D. Sharma, and Z. J. Chong, "Amvnet: Assertion-based multi-view fusion network for lidar semantic segmentation," *arXiv preprint arXiv:2012.04934*, 2020.
- [32] J. Xu, R. Zhang, J. Dou, Y. Zhu, J. Sun, and S. Pu, "Rpnvnet: A deep and efficient range-point-voxel fusion network for lidar point cloud segmentation," 2021.
- [33] H. Caesar, V. Bankiti, A. Lang, S. Vora, V. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nusenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [34] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [35] T. Cortinhal, G. Tzelepis, and E. E. Aksoy, "Salsanext: fast, uncertainty-aware semantic segmentation of lidar point clouds for autonomous driving," *arXiv preprint arXiv:2003.03653*, 2020.
- [36] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 945–953.
- [37] C. Siheng, T. Dong, F. Chen, V. Anthony, and K. Jelena, "Fast resampling of three-dimensional point clouds via graphs," *IEEE Transactions on Signal Processing*, vol. 66, pp. 666–681, 2018.
- [38] J. Peng and C.-C. J. Kuo, "Geometry-guided progressive lossless 3d mesh coding with octree (ot) decomposition," in *ACM SIGGRAPH 2005 Papers*, ser. SIGGRAPH '05, 2005, p. 609–616.
- [39] Z. Songyang, C. Shuguang, and D. Zhi, "Hypergraph spectral analysis and processing in 3d point cloud," *IEEE Transactions on Image Processing*, vol. 30, pp. 1193–1206, 2021.