

3DGTN: 3D Dual-Attention GLocal Transformer Network for Point Cloud Classification and Segmentation

Dening Lu, Kyle (Yilin) Gao, *Graduate Student Member, IEEE*, Qian Xie, Linlin Xu, *Member, IEEE*, Jonathan Li, *Fellow, IEEE*

Abstract—Although the application of Transformers to 3D point cloud processing has achieved significant progress and success, it is still challenging for existing 3D Transformer methods to efficiently and accurately learn both valuable global and local features for improved applications. This paper presents a novel point cloud representational learning network, called 3D Dual Self-attention Global Local (GLocal) Transformer Network (3DGTN), for improved feature learning in both classification and segmentation tasks, with the following key contributions. First, a GLocal Feature Learning (GFL) block with the dual self-attention mechanism (i.e., a novel Point-Patch Self-Attention, called PPSA, and a channel-wise self-attention) is designed to efficiently learn the global and local context information. Secondly, the GFL block is integrated with a multi-scale Graph Convolution-based Local Feature Aggregation (LFA) block, leading to a Global-Local (GLocal) information extraction module that can efficiently capture critical information. Third, a series of GLocal modules are used to construct a new hierarchical encoder-decoder structure to enable the learning of information in different scales in a hierarchical manner. The proposed framework is evaluated on both classification and segmentation datasets, demonstrating that the proposed method is capable of outperforming many state-of-the-art methods on both synthetic and LiDAR data. *Our code has been released at <https://github.com/d62lu/3DGTN>.*

Index Terms—Transformer, Graph convolution, Point cloud classification, Point cloud segmentation, LiDAR data processing, Self-attention mechanism.

I. INTRODUCTION

PPOINT cloud classification and segmentation are fundamental tasks in 3D computer vision. Point clouds, being flexible, simple, and with easy-to-use data structures, are commonly used in 3D mapping, robotics, autonomous navigation, and city information modeling. From the perspective of point cloud processing, both local and global features play an important role in classification and segmentation tasks. Local features refer to the features that capture the local geometric

patterns and details of the point cloud. Global features refer to the features that capture the overall shape and structure of the entire point cloud. A combination of global and local features (called *GLocal* features here) is able to provide the model with a more complete representation of the target point cloud.

For classification and segmentation tasks, many types of deep learning architectures have been experimented with in the recent past. Among these, the Transformer [1] architecture emerged as a powerful point cloud feature extraction backbone, performing exceedingly well on LiDAR point cloud classification and segmentation. [2]–[5]. First developed for natural language processing, the Transformer is a low-inductive bias network that is capable of learning long-range features. Since then, Transformers have successfully been applied to 2D and 3D computer vision to various tasks, achieving state-of-the-art results across a wide variety of benchmarks.

Although existing 3D Transformer approaches demonstrated strong feature learning capabilities in 3D point cloud applications, they still have limitations in terms of modeling both the local information and global information in an efficient and accurate manner. This paper presents a 3D Dual-attention Global-Local (GLocal) Transformer Network, called 3DGTN. It focuses on addressing the difficulty of effectively exploiting global and local features for point cloud classification and segmentation. Many current Transformer methods either emphasize local information extraction or struggle to integrate global and local features accurately. The proposed PPSA mechanism in 3DGTN aims to overcome this limitation. 3DGTN is tailor-designed to improve combined global and local feature learning in 3D point cloud data processing, with the following key characteristics.

- A GLocal Feature Learning (GFL) block with the dual self-attention mechanism is designed to efficiently learn the GLocal context information. The PPSA approach can better capture global correlation among local neighborhoods. The dual-attention mechanism integrates PPSA and Channel-wise Self-Attention (CSA) to improve the learning of critical information in both the spatial domain and feature domain.
- The GFL block is integrated with a Local Feature Aggregation (LFA) block into a Global-Local (GLocal) information extraction module to enable the learning of both valuable global information and critical local information. The LFA block is designed based on the Graph Convolution Network (GCN) to improve both the

This work was supported in part by the Natural Sciences and Engineering Research Council of Canada (NSERC) under Grant RGPIN-2019-06744. The first author was funded by the China Scholarship Council under Ph.D. Scholarship 202106830030. (*Corresponding authors: Linlin Xu*(l44xu@uwaterloo.ca); *Jonathan Li*(junli@uwaterloo.ca).

Dening Lu, Kyle Gao, Linlin Xu, and Jonathan Li are with the Department of Systems Design Engineering, University of Waterloo, Waterloo, Ontario, N2L 3G1, Canada (e-mail: d62lu, y56gao, l44xu, junli@uwaterloo.ca).

Jonathan Li is also with the Department of Geography and Environmental Management, University of Waterloo, Waterloo, Ontario, N2L 3G1, Canada.

Qian Xie is with the Department of Computer Science, University of Oxford, Oxford, OX1 3QD, U.K. (e-mail: qian.xie@cs.ox.ac.uk).

efficiency and accuracy of local information extraction.

- The GLocal modules are used to construct a new hierarchical encoder-decoder structure to enable the learning of information at different scales in a hierarchical manner, leading to a general point cloud representation network that can improve both classification and segmentation.

Extensive experiments comparing the proposed approach with many state-of-the-art algorithms on many datasets, i.e., ModelNet40, ScanObjectNN, ShapeNet, and Titan MultiSpectral (MS) LiDAR datasets, demonstrate that our method exceeds previous state-of-the-art performance in both classification and segmentation tasks.

II. RELATED WORK

Transformer-based methods tailored for point cloud data can be broadly categorized into two main groups: global Transformer-based methods and local Transformer-based methods. Here, we review existing approaches in both categories and summarize the limitations.

A. Global Transformers in 3D Point Cloud Processing

The global Transformer approaches focus on learning large-scale context information from the 3D point cloud to improve classification and segmentation. Point Cloud Transformer (PCT), as a standard global Transformer network, was proposed in [6]. In PCT, all input points were leveraged for global feature extraction. PCT first adopted a neighborhood-embedding strategy to aggregate the local information, followed by feeding the embedded features into four stacked global Transformer blocks. Lastly, it utilized a global Max and Average (MA) pooling to extract the global information for point cloud classification. The segmentation network variant of PCT [6] had the same feature encoding backbone as the classification network variant. However, the decoder first concatenated the pooled global feature with each point feature, enhancing the perception of global information for each point. Then the concatenated features were fed into a series of MLP layers for dense prediction, following PointNet [7].

3CROSSNet proposed in [8] used multi-scale global information for classification. Taking the raw point cloud as input, it first generated three point subsets with different resolutions using Farthest Point Sampling (FPS). Secondly, it established k -Nearest Neighborhood (k NN) [7] and extracted local information using a series of Multi-Layer Perception (MLP) modules for each point subset. Thirdly, the cascaded global Transformer blocks were applied to extract the global information of each subset. Lastly, given the multi-scale global features, 3CROSSNet used the Cross-Level Cross-Attention (CLCA) and Cross-Scale Cross-Attention (CSCA) modules to capture long-range inter- and intra-level dependencies for classification.

Instead of using raw point clouds, Stratified Transformer [9] took 3D voxels as input to the segmentation network. It applied Transformer blocks in predefined local windows, following Swin Transformer [10]. To capture the global information and establish connections between different windows, it presented a novel *key* sampling strategy, enlarging the effective receptive field for each *query* point.

B. Local Transformers in 3D Point Cloud Processing

As a local Transformer network, Point Transformer (PT) [11] focused on extracting local information by the Transformer. A downsampled pointset was passed through five local Transformer blocks. Specifically, for each block, PT used k NN for sampling points, then utilized a vector-attention mechanism to capture local features. After five local Transformer blocks, PT used a global MA pooling to extract the global feature for classification. Local Feature Transformer Network (LFT-Net) [12] had a similar architecture. However, it used an additional trans-pooling module to alleviate the feature loss during the pooling. For 3D point cloud segmentation, PT [11] developed the segmentation network based on its classification framework. The authors designed a U-net-style architecture for segmentation, where the decoder was symmetric to the encoder. Since it used a hierarchical structure in the encoder, a transition-up module with trilinear interpolation was proposed in the decoder for point cloud upsampling.

C. Limitations of Current 3D Transformer-based Networks

Despite the great success of Transformers in point cloud classification and segmentation, existing 3D Transformer methods tend to only consider local information extraction or struggle to learn both global and local features effectively. This issue makes it still challenging for 3D Transformer methods to capture the global information of the target accurately while preserving the local features. For example, PT [11] only utilized Transformer blocks in local neighborhoods, while ignoring global feature learning. FlatFormer [13] used Transformer blocks to extract window-based local features, and designed a window shift strategy to indirectly achieve global feature learning. PCT [6] only captured local information at the beginning of the network, as data preprocessing. It cannot dynamically fuse the local information with the global information extracted from each stage in the network. Recently, there have been several works [14]–[19] that extract both local and global features in a simple cascading way. PatchFormer [17], SPFormer [18], and SPT [19] all used Transformer blocks to capture global features from aggregated superpoint-based local features. However, it is easy for them to lose local neighborhood information. Therefore, this paper proposes a novel PPSA (Sec. III-C) mechanism to improve global and local feature learning. It aims to explicitly fuse the local neighborhood and global information of the target. To our knowledge, our 3DGTN is the first work to introduce combined global-local feature learning to 3D Transformers for point cloud processing.

III. 3D DUAL-ATTENTION GLOCAL TRANSFORMER

In this section, we introduce the encoder and decoder structures of our 3DGTN for both point cloud classification and segmentation. We first show the pipeline of our method, then introduce the main blocks in the encoder and decoder respectively.

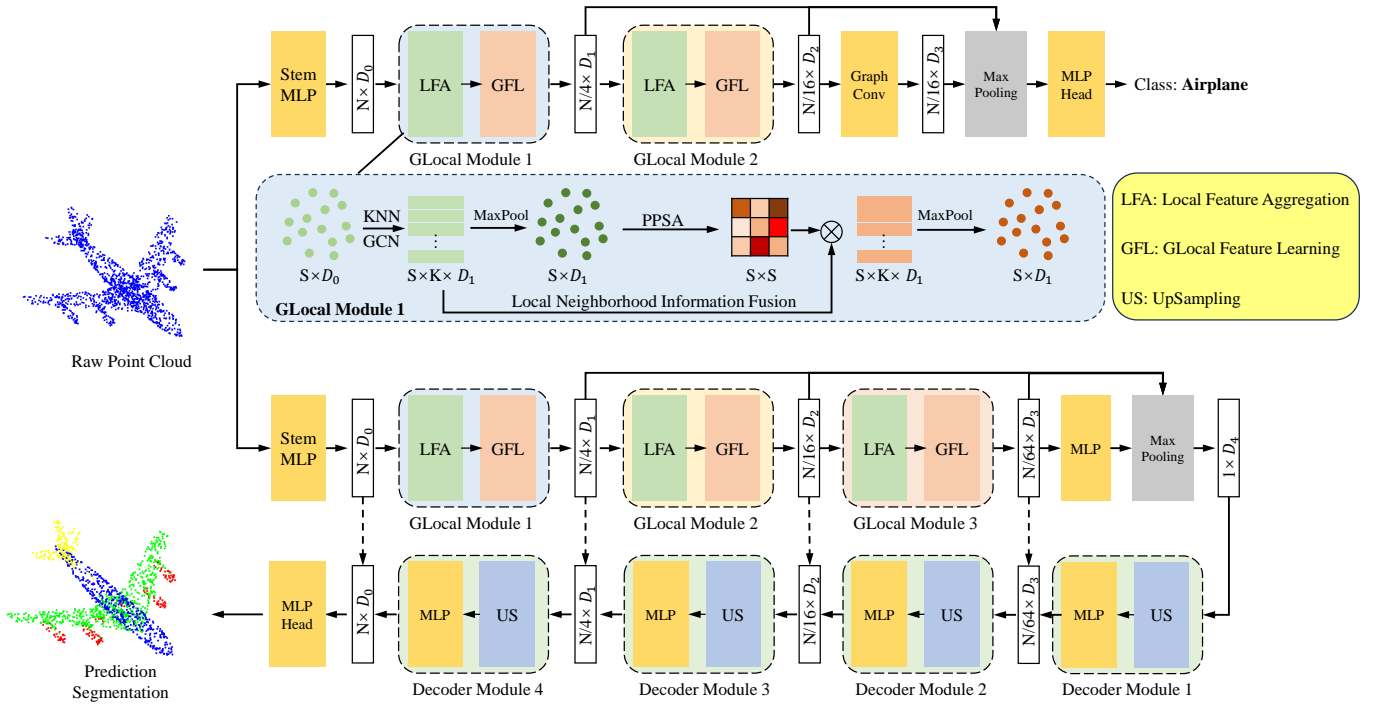


Fig. 1. 3DGTN networks for point cloud classification (top) and segmentation (bottom), where GCN-based LFA blocks and dual self-attention-based GFL blocks are designed for strong feature representation. A brief illustration of the GLocal Module 1 is provided to highlight GLocal feature generation. Please refer to Fig. 2, Sec. III-B and III-C for more details.

A. Overview

Fig. 1 shows the overall pipeline of our method. Our classification and segmentation networks use the same encoder architecture. After that, the classification network utilizes an MLP head to obtain the final classification results, while the segmentation network utilizes a decoder with trilinear interpolation-based upsampling for dense prediction.

The original point cloud is taken as input to the encoder. We first design a stem MLP block to project the input data into a higher-dimensional space. After that, the projected features are fed into stacked LFA and GFL blocks in a hierarchical manner for GLocal feature extraction. Specifically, the LFA block is adapted from the multi-scale GCN [20], and the GFL block is adapted from the Transformer. Following this, we use the max-pooling operation on the output feature maps of each module, to obtain the GLocal feature of each level. Then, we concatenate them for multi-level GLocal feature generation. Given the extracted feature, we leverage an MLP head for the point cloud classification task, which consists of two fully connected layers with batch normalization and RELU activation. For the segmentation task, the extracted features are then taken as input to the decoder. To improve efficiency, we adopt an ALL-MLP decoder structure, instead of a symmetric one. In the upsampling block, the interpolated points are concatenated with the corresponding feature points from the encoder via a skip connection. The trilinear upsampling method we used has been widely applied to hierarchical networks of point cloud processing. It generates new points by considering the weighted averages of neighboring points

in the geometric space as 3D linear interpolation, providing an effective method to enhance the density and precision of 3D data representations. We note that the number of modules in the encoder and decoder can vary according to the number of input points. In our experiments, we designed a two-module encoder for the classification task (1024 points), but a three-module encoder and corresponding decoder for the segmentation task (2048 points).

B. Local Feature Aggregation Block

We adopt the GCN-based LFA block for local feature aggregation. The LFA block (Fig. 2) is introduced as follows.

The input point cloud is first downsampled to $N/4$ points via FPS, generating a sampled point subset S , where N is the number of the input points. After that, the LFA block constructs multi-scale k -NN neighborhoods (three scales k_1, k_2, k_3 in our experiments) for each sampled point, to ensure the diversity of the receptive fields. In each neighborhood χ_i of the sampled point S_i , a fused feature \mathbb{C}_{ij} is generated by computing the difference between the j -th neighborhood point χ_{ij} in χ_i and S_i :

$$\mathbb{C}_{ij} = \text{concat}(\mathbb{F}_{ij} - \mathbb{F}_i, \mathbb{F}_i), \quad (1)$$

where \mathbb{F}_i and \mathbb{F}_{ij} represents the feature of S_i and χ_{ij} respectively. Given the fused neighborhood feature, the Graph Convolution in χ_i can be formulated as:

$$l_i = \text{maxpooling}(\text{Conv}(\mathbb{C}_{ij})), \quad (2)$$

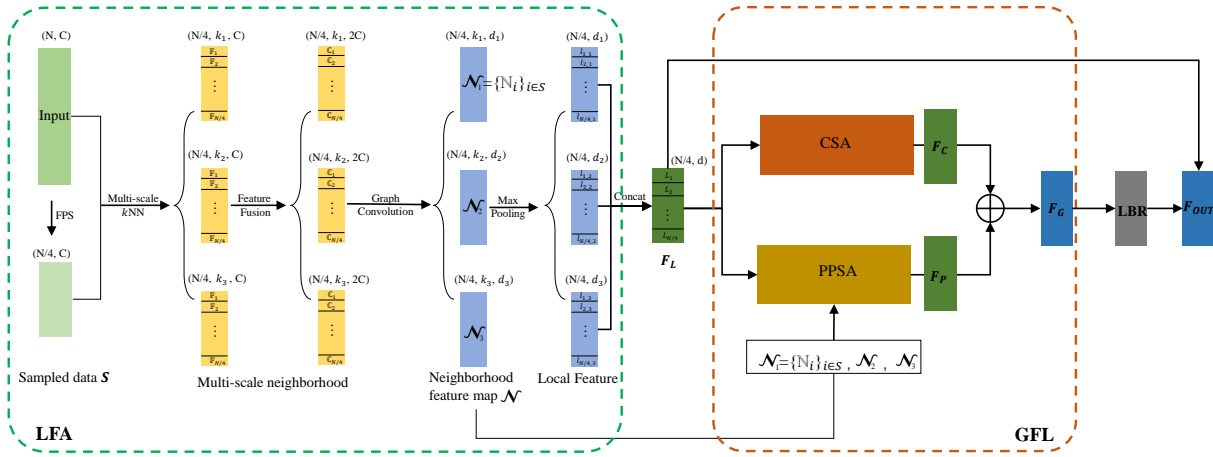


Fig. 2. Architecture of GLocal Module 1, which consists of an LFA block and a GFL block.

where l_i is the aggregated local feature of S_i , $Conv$ is a convolution operation with 1×1 kernels. Specifically, in Fig. 2, we denote the dimension of the input feature map as (N, C) , and the dimension of \mathbb{C}_{ij} is $2C$. Furthermore, we define different output dimensions of Graph Convolution for different scale neighborhoods: d_1, d_2 , and d_3 , where $d_1 < d_2 < d_3$ ($k_1 < k_2 < k_3$). $Conv(\mathbb{C}_{ij})$ establishes semantic relationships between the sampling point S_i and neighborhood point χ_{ij} . As such, a neighborhood feature set containing local information, \mathbb{N}_i of S_i , is generated. Then, the max-pooling operation is used to aggregate the local information to S_i .

The multi-scale local feature L_i of S_i , can be expressed as via a concatenation as:

$$L_i = \text{concat}(l_{i_1}, l_{i_2}, l_{i_3}), \quad (3)$$

where $l_{i_1}, l_{i_2}, l_{i_3}$ represent three local features of S_i at three different scales.

C. GLocal Feature Learning Block

Our GFL block contains two kinds of self-attention mechanisms: PPSA and CSA. PPSA, as a novel point-wise self-attention mechanism, is proposed to fuse the global features and local neighborhood information extracted from the LFA block for better GLocal feature learning. CSA is utilized to measure the correlation among different feature channels. It is able to improve context information modeling by highlighting the role of interaction across various channels. A detailed introduction to these two mechanisms is as follows.

Point-Patch Self-Attention. PPSA fuses local and global features. As shown in Fig. 3, the aggregated features $F_L = \{L_i\}_{i \in S} \in R^{s \times d}$ from the LFA block is taken as input, where s is the number of sampled points in S , and d denotes the feature dimension of F_L . We first project F_L into two different feature spaces to generate *Query*, *Key* matrices:

$$\begin{aligned} \text{Query} &= F_L W_{QP}, \\ \text{Key} &= F_L W_{KP}, \end{aligned} \quad (4)$$

where W_{QP}, W_{KP} are learnable weight matrices. Then, the attention map $M_P \in R^{s \times s}$ of PPSA can be formulated as:

$$M_P = \text{softmax}\left(\frac{QK^T}{\sqrt{d}} + B\right), \quad (5)$$

where Q, K denote the *Query*, *Key* matrices, and B is a learnable position encoding matrix defined by [11]. Next, we treat the neighborhood feature map $\mathcal{N} = \{\mathbb{N}_i\}_{i \in S}$ at each scale as the *Value* branch, instead of F_L used by the vanilla PSA. In other words, the elements in the attention map are taken as weights of the corresponding neighborhood feature sets in \mathcal{N} . Then, the output neighborhood feature set is obtained by computing a weighted sum of all input sets. As such, we leverage all the points including sampling points and neighborhood points for the GLocal information extraction, instead of only sampling points. This method is able to improve the feature learning and mitigate the local information loss caused by the pooling operation in Eq. 2. Given the aforementioned attention map M_P and the *Value* matrix, the output GLocal feature can be expressed as:

$$F_o = \text{maxpooling}(M_P V), \quad (6)$$

where V denotes the *Value* matrix, i.e., \mathcal{N} . The detailed algorithm flow and feature dimension transformation of PPSA are shown in Fig. 3. We note that there are three neighborhood feature maps \mathcal{N} for each sampled point S_i because of the multi-scale grouping strategy, which are denoted as $\mathcal{N}_1, \mathcal{N}_2$, and \mathcal{N}_3 . Correspondingly, we obtain three output GLocal features at different scales, F_{o1}, F_{o2} , and F_{o3} . Lastly, we concatenate them to get the final point-wise GLocal feature F_P :

$$F_P = \text{concat}(F_{o1}, F_{o2}, F_{o3}), \quad (7)$$

where F_{o1}, F_{o2}, F_{o3} are generated from neighborhood feature maps at different scales.

Channel-wise Self-Attention. Apart from the PPSA mechanism, we also utilize the CSA mechanism to capture context dependencies in the channel dimension. It enables the model to build connections among different channels, allowing it to focus on different feature channels depending on input

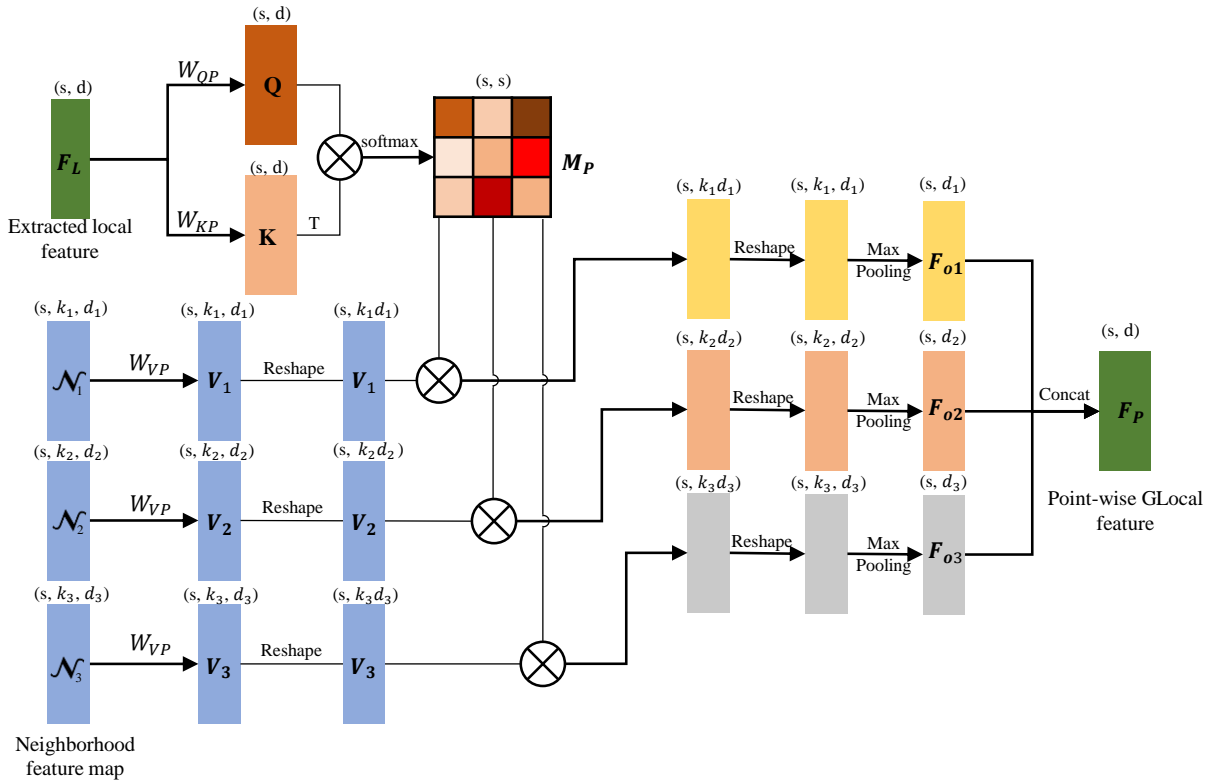


Fig. 3. Point-Patch Self-Attention mechanism. It utilizes the sampling point features and the corresponding neighborhood (patch) feature maps for point-wise GLocal feature extraction.

data. As shown in Fig. 4, given the aggregated local feature $F_L \in R^{s \times d}$, we first compute the attention map $M_C \in R^{d \times d}$ of CSA as:

$$M_C = K^T Q = (F_L W_{KC})^T (F_L W_{QC}). \quad (8)$$

where the shapes of K, Q are reduced to $(s/8) \times d$ by weight matrices W_{KC} and W_{QC} , to improve efficiency. Inspired by [21], we calculate the affinity matrix A_C based on M_C , to measure the difference among channels, which can be expressed as:

$$A_C = \text{softmax}(\text{expand}(\text{maxpooling}(M_C)) - M_C), \quad (9)$$

where $\text{maxpooling}(M_C) \in R^{d \times 1}$ extracts the maximum value of each row in M_C , $\text{expand}(\cdot)$ expands the matrix $\text{maxpooling}(M_C)$ to the same size as M_C by column repetition. From the subtraction, the larger in magnitude the element in A_C , the lower the similarity of the corresponding two channels. As such, CSA tends to focus on channels with significant differences, avoiding aggregating similar/redundant information. After that, we calculate the *Value* matrix as:

$$V = F_L W_{VC}, \quad (10)$$

where W_{VC} is a learnable weight matrix. Finally, the channel-wise GLocal feature F_C can be expressed as:

$$F_C = V A_C. \quad (11)$$

Given both F_P and F_C , the final GLocal feature can be generated by combining them with an element-wise addition:

$$F_G = F_P + F_C. \quad (12)$$

Additionally, we apply a residual connection between the LFA block and the GFL block:

$$F_{OUT} = F_L + LBR(F_G), \quad (13)$$

where F_{OUT} is the final output feature map of the defined GFL block, and LBR denotes the combination of *Linear*, *BatchNorm*, and *ReLU* layers.

IV. EXPERIMENTS

In this section, we first introduce the implementation details of our 3DGTN, including hardware configuration and hyperparameter settings. Secondly, we present the performance evaluation of our method on classification and segmentation tasks, comparing it to state-of-the-art methods. Specifically, we tested our method for the classification task on the widely-used ModelNet40 and ScanObjectNN datasets [22], [43]. For object part segmentation, we tested our method on the ShapeNet dataset [44]. For semantic segmentation, we tested our method on the challenging large-scale MS-LiDAR dataset [45]. Finally, we present the ablation experiment results on the main components of our method.

A. Implementation Details

We implemented our classification and segmentation networks in PyTorch. Both were trained and tested on an NVIDIA Tesla V100 GPU. We used the SGD Optimizer with a momentum of 0.9 and weight decay of 0.0001. The initial learning rate was set to 0.01, with a cosine annealing schedule. We trained

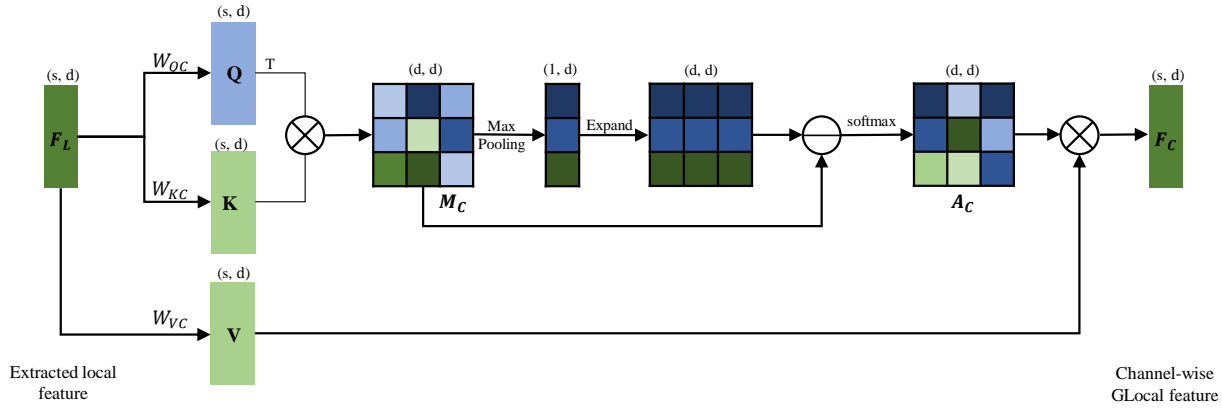


Fig. 4. Channel-wise self-attention mechanism. An Affinity matrix A_C is designed to avoid aggregating redundant features, enhancing the channel-wise GLocal feature representation.

TABLE I
QUANTITATIVE COMPARISON (mAcc, OA, PARAMETERS, FLOPS, AND FRAME PER SECOND) OF CLASSIFICATION PERFORMANCE ON THE MODELNET40 DATASET. TRANSFORMER METHODS ARE SEPARATED FROM OTHER LEARNING-BASED METHODS. THE HIGHEST EVALUATION SCORE IS SHOWN IN BOLD TYPE.

Methods	Input Size	mAcc (%)	OA (%)	Parameters (MB)	FLOPs (GB)	Frame Per Sec.
Other Learning-based Methods						
3DShapeNets [22]	1024	77.3	84.7	-	-	-
PointNet [7]	1024	86.0	89.2	3.47	0.45	614
PointNet++ [23]	1024	88.2	91.9	1.74	4.09	16
diffConv [24]	1024	90.4	93.2	2.08	0.16	-
CurveNet [25]	1024	90.4	93.1	-	-	-
PointCNN [26]	1024	88.1	92.2	0.6	1.54	14
DGCNN [20]	1024	90.2	92.2	1.81	2.43	279
FatNet [27]	1024	90.6	93.2	-	-	-
DRNet [28]	1024	-	93.1	-	-	-
PointMLP [29]	1024	91.4	94.1	12.6	-	112
Point-PN [30]	1024	-	93.8	0.8	-	532
RepSurf [31]	1024	91.1	94.0	1.48	-	205
Transformer-based Methods						
PATs [32]	1024	-	91.7	-	-	-
LFT-Net [12]	1024	89.7	93.2	-	-	-
PointTransformer [33]	1024	89.0	92.8	13.86	9.36	17
MLMST [34]	1024	-	92.9	-	-	-
PointCloudTransformer [6]	1024	90.3	93.2	2.80	2.02	125
LSLPCT [35]	1024	90.5	93.5	-	-	-
PointTransformer [11]	1024	90.6	93.7	9.14	17.14	15
CloudTransformers [36]	1024	90.8	93.1	22.91	12.69	12
GBNet [21]	1024	91.0	93.8	8.38	9.02	102
3DCTN [14]	1024	91.6	93.2	4.21	3.76	24
PatchFormer [17]	1024	-	93.5	2.45	1.62	201
Ours	1024	92.4	94.0	5.21	3.09	15

classification, part segmentation, and semantic segmentation networks for 250, 300, and 500 epochs respectively, with the same batch size of 16.

B. Point Cloud Classification

Datasets and Metrics. The ModelNet40 dataset contains 12311 CAD models with 40 object categories. We split them into 9843 training samples and 2468 testing models, following PointNet++ [23]. For a fair comparison, we downsampled each input point cloud to 1024 points with normals via FPS. Since point clouds in ModelNet40 are generated from 3D meshes, we can easily obtain the normal of each point according to the corresponding surface normal. The mean accuracy

within each category (mAcc) and the overall accuracy (OA) are used for performance evaluation, which are formulated as:

$$mAcc = \frac{\sum_{i=1}^K \frac{T_i}{N_i}}{K}, \quad (14)$$

$$OA = \frac{T}{N},$$

where T is the number of all correctly predicted point clouds, $T = \sum_{i=1}^K T_i$, T_i is the number of correctly predicted point clouds in class i , K is the number of classes in the dataset, N is the number of all point clouds in the dataset, $N = \sum_{i=1}^K N_i$ and N_i is the number of point clouds in class i . Additionally, we adopt the total number of parameters, FLOPs (FLOating Point operations), and Frame Per Second to evaluate the model size and efficiency.

TABLE II
QUANTITATIVE COMPARISON (%) OF CLASSIFICATION PERFORMANCE ON THE SCANOBJECTNN DATASET. TRANSFORMER-BASED METHODS AND OTHER LEARNING-BASED METHODS ARE SEPARATED. THE HIGHEST EVALUATION SCORE IS SHOWN IN BOLD TYPE.

Methods	Input Size	mAcc (%)	OA (%)
Other Learning-based Methods			
PointNet [7]	1024	63.4	68.2
PointNet++ [23]	1024	75.4	77.9
SpiderCNN [37]	1024	69.8	73.7
PointCNN [26]	1024	75.1	78.5
DGCNN [20]	1024	73.6	78.1
SimpleView [38]	1024	-	80.5
PRANet [39]	1024	79.1	82.1
PointMLP [29]	1024	84.4	85.7
RepSurf [31]	1024	83.1	86.0
Transformer-based Methods			
PointTransformer [33]	1024	75.3	77.6
PointCloudTransformer [6]	1024	77.1	80.5
PointTransformer [11]	1024	78.2	80.8
GBNet [21]	1024	77.8	80.5
Point-MAE [40]	1024	-	85.2
Point-TnT [41]	1024	-	83.5
Point-BERT [42]	1024	-	83.1
3DCTN [14]	1024	79.5	81.5
Ours	1024	83.2	85.8

To further evaluate the performance of 3DGTN to the real-world data captured by LiDAR scanning, ScanObjectNN [43] classification performance was also tested in our experiments. There are $\sim 15,000$ objects in ScanObjectNN, which are categorized into 15 categories with 2902 unique object instances. Since each object was segmented from the scanned scene point cloud, object point clouds usually include numerous outliers in the form of background points, and were corrupted by occlusions and noises. Therefore, it was more challenging to perform shape classification on this dataset. We used the hardest variant of the dataset (*PB_T50_RS*) and adopted the original training/testing split as in [43]. Similarly, each sample from ScanObjectNN was downsampled to 1024 points. Since point clouds in *PB_T50_RS* have no normal information, we only took the 3D coordinates of point clouds as input.

Performance Comparison. We compared our 3DGTN with the state-of-the-art Transformer-based methods and other

deep learning-based methods. The comparison results are shown in Table. I and II. Specifically, for the ModelNet40 dataset, our method achieves the best mean accuracy of 92.4% among all benchmarked methods in terms of mAcc, outperforming the prior state-of-the-art (PointMLP [29]) by 1.0 absolute percentage points. In terms of OA, our method achieves the best result of 94.0% among the Transformer-based methods. For the model size, our method requires fewer parameters (5.21MB) and FLOPs (3.09GB) compared to most Transformer-based algorithms, accounting for only 57% and 18% of Point Transformer [11] respectively. However, due to the naive implementation of several time-consuming operations like downsampling and k NN neighborhood construction, the inference speed of our method can still be improved. For the ScanobjectNN dataset, 3DGTN also achieves competitive performance with the SOTA approaches. Especially, it obtains the best results in terms of both OA (85.8%) and mAcc (83.2%) among all compared Transformer-based methods, which demonstrates the excellent performance of 3DGTN in LiDAR data processing.

Visualization. We generate Grad-CAM [46] map visualization results from the ModelNet40 dataset. The Grad-CAM technique is designed to produce a coarse localization map highlighting the important regions in target point clouds. It uses the gradient information flowing into the last convolutional layer of a deep network to understand the importance of each neuron for a decision of interest. As shown in Fig. 5, we obtain the regions of interest of our network for several point clouds of the Airplane, Car, Cup, and Plant classes. From the results, the attention (colored in red) is mainly focused on the wings and tail of the Airplane, the tires of the Car, the handle of the Cup, and the leaves of the Plant. As we can see, all the regions of interest are consistent with the human visual system, which helps us establish appropriate trust in predictions from deep networks.

C. Part Segmentation on ShapeNet Dataset

Dataset and Metrics. The ShapeNet dataset contains 16880 models with 16 shape categories. We split them into 14006

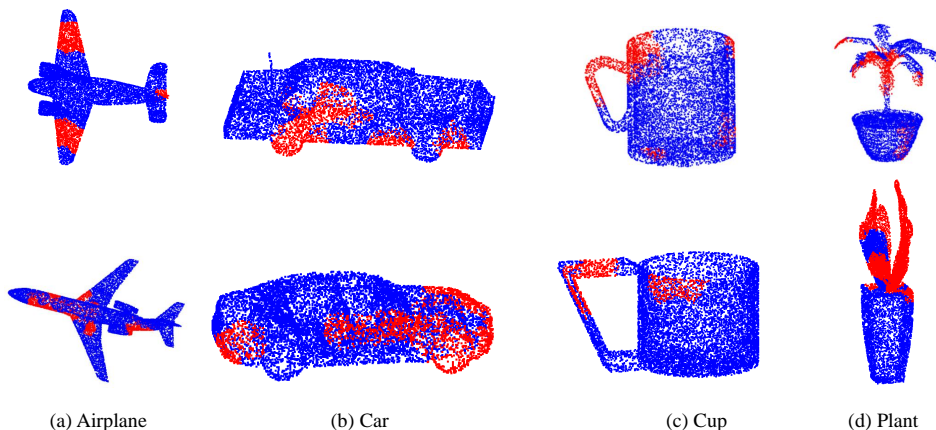


Fig. 5. Visualization of 3DGTN attention on the ModelNet40 classification dataset. As can be seen, the attention (red) is focused on the discriminative parts of targets, such as the wings of an airplane, the tires of a car, the handle of a cup, and the leaves of a plant.

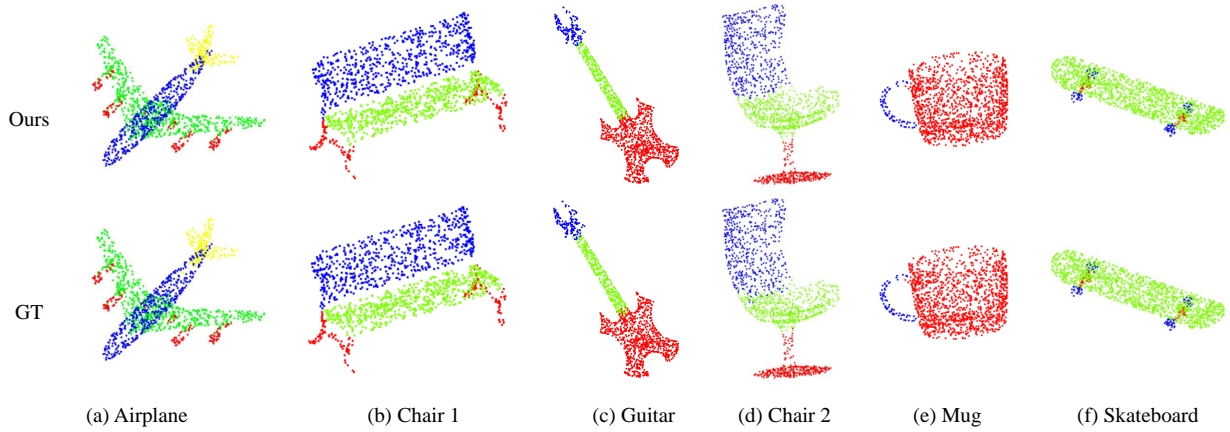


Fig. 6. Part segmentation results from the ShapeNet dataset. As can be seen, our segmentation predictions are faithful to the ground truth.

TABLE III
QUANTITATIVE COMPARISON (%) OF PART SEGMENTATION PERFORMANCE ON THE SHAPENET DATASET. THE HIGHEST EVALUATION SCORE IS SHOWN IN BOLD.

Methods	Cat. mIoU	Ins. mIoU
PointNet [7]	80.4	83.7
PointNet++ [23]	81.9	85.1
A-SCN [47]	-	84.6
PCNN [48]	81.8	85.1
SpiderCNN [37]	82.4	85.3
SPLATNet [49]	83.7	85.4
SGPN [50]	82.8	85.8
SubSparseCNN [51]	83.3	86.0
PointCNN [26]	84.6	86.1
PointConv [52]	82.8	85.7
Point2Sequence [53]	-	85.2
DGCNN [20]	82.3	85.2
PVCNN [54]	-	86.2
RS-CNN [55]	84.0	86.2
KPConv [56]	85.0	86.2
InterpCNN [57]	84.0	86.3
DensePoint [58]	84.2	86.4
PACConv [59]	84.6	86.1
PointTransformer [11]	83.7	86.6
StratifiedTransformer [9]	85.1	86.6
PatchFormer [17]	-	86.5
APES [60]	83.7	85.8
Ours	84.0	86.6

training samples and 2874 testing models, following Point Transformer [11]. The dataset has 50 part labels, and each object has at least two parts. For a fair comparison, we downsampled each input point cloud to 2048 points with normals by FPS. The category-wise mean Intersection over Union (mIoU) and instance-wise mIoU [11] are used for performance evaluation, which are formulated as below:

$$\begin{aligned}
 Cat.mIoU &= \frac{\sum_{i=1}^{Cls} \sum_{j=1}^{H_i} mIoU_j}{Cls}, \\
 mIoU_j &= \frac{\sum_{i=1}^{M_j} \frac{TP_i}{TP_i + FP_i + FN_i}}{M}, \\
 Ins.mIoU &= \frac{\sum_{i=1}^G mIoU_i}{G},
 \end{aligned} \quad (15)$$

where Cls is the number of total shape classes of the dataset ($Cls = 16$ in the ShapeNet dataset), H_i represents the number

of instances of the class i , M_j represents the number of part classes (varies with shape classes) in the j -th instance, TP_i represents the number of the true positive samples in the i -th part class, and G is the numbers of all instances in the dataset ($G = 2874$ in the testing dataset of ShapeNet.)

Performance Comparison. The comparison results are shown in Table. III. As measured by instance-wise mIoU, our 3DGTN achieves competitive results (86.6%) compared with the SOTA Transformer-based methods such as Stratified Transformer [9]. This demonstrates the excellent performance of 3DGTN in terms of part segmentation. Several part segmentation results are shown in Fig. 6.

D. Semantic Segmentation on Airborne MS-LiDAR Dataset

Dataset and Metrics. Most recently, a large-scale airborne MS-LiDAR dataset was proposed in [45]. We tested 3DGTN on this dataset to explore its performance in practical remote sensing applications. The MS-LiDAR dataset was captured by a Teledyne Optech Titan MS-LiDAR system [45]. In addition to three-dimensional coordinates, each point also has three channels with wavelengths of 1,550 nm (*MIR*), 1,064 nm (*NIR*), and 532 nm (*Green*). The dataset was labelled into six categories: Road, Building, Grass, Tree, Soil, and Powerline. The dataset was divided into 13 subsets, where subsets 1-10 were taken as training data, while subsets 11-13 were taken as testing data. For fair comparison, we took the same data pre-processing (data fusion, normalization, and training/testing sample generation) methods described in [45]. The average F_1 score [63], mIoU, and OA are used for performance evaluation:

$$\begin{aligned}
 mIoU &= \frac{\sum_{i=1}^{Cls} mIoU_i}{Cls}, \\
 mIoU_i &= \frac{\sum TP_i}{\sum TP_i + \sum FP_i + \sum FN_i}, \\
 AverageF_1 &= \frac{\sum_{i=1}^{Cls} F_{1i}}{Cls},
 \end{aligned} \quad (16)$$

TABLE IV

CONFUSION MATRIX OF 3DGTN ON THE AIRBORNE MS-LIDAR DATASET. THE SECOND TO SEVENTH ROW REPRESENT THE NUMBER OF POINTS, THE LAST THREE ROWS REPRESENT THE PRECISION, RECALL, AND F_1 SCORE IN % FOR EACH CLASS

Categories	Road	Building	Grass	Tree	Soil	Powerline
Road	200355	14692	49	7	6274	0
Building	29709	871560	11144	612	14449	0
Grass	203	7869	925649	1281	18	424
Tree	71	648	2513	108529	2	14
Soil	6691	10628	103	162	26802	0
Powerline	5	15	323	0	0	7036
Precision	90.50	93.97	98.95	97.09	60.38	95.35
Recall	84.53	96.26	98.50	98.14	56.37	94.14
F_1	87.44	95.10	98.73	97.61	58.31	94.74

TABLE V

QUANTITATIVE COMPARISON (%) OF SEMANTIC SEGMENTATION PERFORMANCE ON THE AIRBORNE MS-LIDAR DATASET. THE HIGHEST EVALUATION SCORE IS SHOWN IN BOLD. THE F_1 SCORE FOR EACH CATEGORY IS ALSO PROVIDED.

Methods	Road	Building	Grass	Tree	Soil	Powerline	Average F_1	mIoU	OA
PointNet [7]	50.81	79.20	68.61	75.21	12.73	22.56	51.52	44.28	83.79
PointNet++ [23]	71.08	83.98	93.24	96.45	30.24	57.28	72.05	58.60	90.09
DGCNN [20]	70.42	90.25	93.62	97.93	21.97	55.24	71.57	51.04	91.36
RSCNN [55]	71.18	89.00	91.42	95.63	26.43	70.03	73.90	56.10	90.99
GACNet [61]	64.51	84.21	93.41	96.66	22.77	33.83	67.65	51.04	89.91
SE-PointNet++ [62]	70.32	85.64	94.70	97.05	37.02	70.35	75.84	60.15	91.16
FR-GCNet [45]	82.63	90.81	95.33	98.77	28.72	74.11	78.61	65.78	93.55
Xiao et al. [63]	73.33	90.51	86.30	95.20	59.24	95.60	83.30	79.25	94.04
GCNAS [64]	87.75	98.68	96.00	99.49	50.74	96.12	88.13	82.23	95.19
Ours	87.44	95.10	98.73	97.61	58.31	94.74	88.63	82.05	95.20

where $Cls = 6$ in the airborne MS-LiDAR dataset, and F_{1i} is calculated as follows:

$$\begin{aligned}
 F_{1i} &= \frac{Precision_i * Recall_i}{Precision_i + Recall_i}, \\
 Precision_i &= \frac{\sum TP_i}{\sum TP_i + \sum FP_i}, \\
 Recall_i &= \frac{\sum TP_i}{\sum TP_i + \sum FN_i}.
 \end{aligned}
 \tag{17}$$

Additionally, the F_{1i} score for each category i is also provided.

Performance Comparison. As shown in Table. IV, the semantic segmentation results of Airborne MS-LiDAR data are presented in the form of a confusion matrix. Since we integrate the CSA mechanism into global feature learning, our method is able to handle multispectral LiDAR point cloud segmentation well. Specifically, from the table, the number of samples differs significantly among categories. In this case of extremely imbalanced data, our 3DGTN still achieves excellent F_1 scores

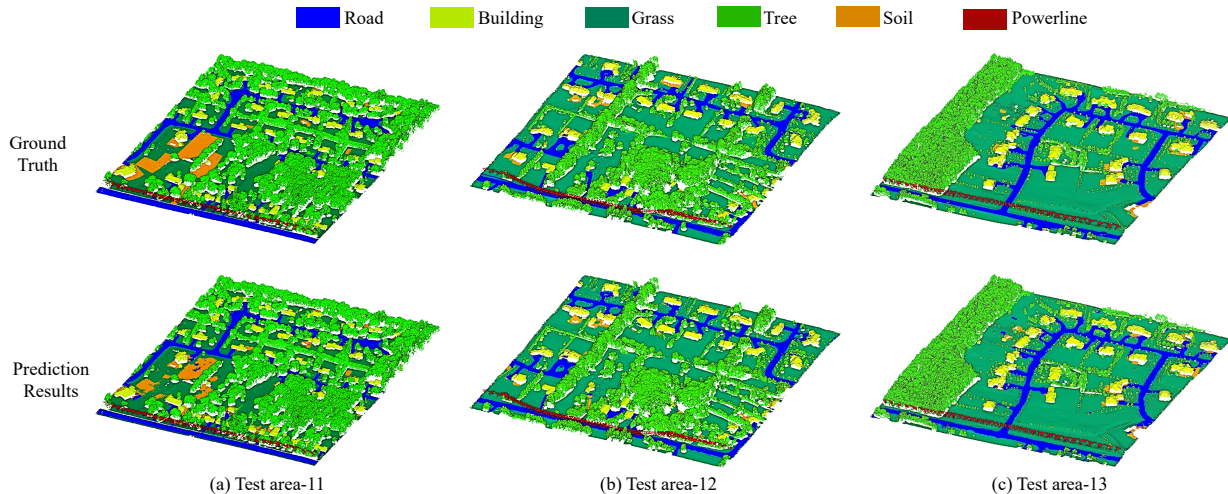


Fig. 7. Semantic segmentation results and ground truth from the Airborne MS-LiDAR dataset.

TABLE VI

RESULTS (%) COMPARISON OF INPUT DATA WITH REMOVAL OF DIFFERENT CHANNEL ON THE AIRBORNE MS-LIDAR DATASET. THE F_1 SCORE FOR EACH CATEGORY IS ALSO PROVIDED

Input				Road	Building	Grass	Tree	Soil	Powerline	Average F_1	mIoU	OA
XYZ	MIR	NIR	Green									
✓	-	✓	✓	84.06	90.70	95.78	96.19	51.22	91.12	84.85	79.92	93.90
✓	✓	-	✓	60.62	84.59	93.30	90.63	33.25	83.37	74.29	69.76	86.41
✓	✓	✓	-	83.08	90.83	92.66	96.95	45.19	88.05	82.79	78.73	92.09
✓	✓	✓	✓	87.44	95.10	98.73	97.61	58.31	94.74	88.63	82.05	95.20

TABLE VII

QUANTITATIVE COMPARISON RESULTS OF ABLATION STUDIES FOR THE MAIN COMPONENTS OF 3DGTN, WHICH WERE PERFORMED ON THE MODELNET40 CLASSIFICATION DATASET. — MEANS COMPONENT REMOVAL, AND → MEANS COMPONENT CHANGING.

Ablation		mAcc (%)	OA (%)	Parameters (MB)	FLOPs (GB)	Frame Per Sec.
Local feature aggregation	Graph convolution → Standard MLP	91.6	93.1	5.18	0.45	17
	Multi-scale → Single-scale	91.3	92.9	2.11	0.64	19
	FPS → Random sampling	92.1	93.5	5.21	3.09	23
GLocal feature learning	—	91.2	92.7	2.25	2.78	17
	— CSA	91.9	93.7	4.71	3.04	16
	— PPSA	91.3	93.2	3.75	2.95	16
	PPSA → Vanilla PSA	91.5	93.6	5.11	3.08	15
Multi-level GLocal feature	—	92.1	93.6	4.75	3.09	15
3DGTN		92.4	94.0	5.21	3.09	15

of over 85% for all categories except soil. The F_1 scores of the grass, tree, and building are over 95%. However, since the geometric characteristics of the soil are very similar to those of grass, which tends to confuse the network, the segmentation results of the soil are not very satisfactory. More feature discrimination approaches would be designed in our future work which could improve the segmentation of similar classes. The comparison results are shown in Table. V. As can be seen, our 3DGTN outperforms all benchmarked methods in terms of average F_1 score (88.63%). It surpasses the prior SOTA methods such as [63] and [64] by 5.33 and 0.50 absolute percentage points, respectively. It also achieves the best OA (95.20%) and a competitive mIoU (82.05%). The prediction results and corresponding ground truth of testing data are shown in Fig. 7. These results demonstrate that our method has excellent performance in processing real-scanned data, exceeding previous SOTA.

We also explored the importance of different channels in the MS-LiDAR data. Specifically, we removed the each of the three channels (*MIR*, *NIR*, *Green*) of MS-LiDAR data, and then analyzed the corresponding performance changes in Table. VI. When the *NIR* channel was removed (Row 4), the performance dropped significantly (average F_1 score was reduced to 74.29% from 88.63%). There is also a slight performance drop when the *MIR* or *Green* channel is removed. The results demonstrate that all these three channels are useful for data segmentation, where the *NIR* channel contributes the most to performance.

E. Ablation Study

We conducted a series of ablation experiments for the main components of our 3DGTN to verify their effectiveness. These experiments were performed on the ModelNet40 dataset.

Local Feature Aggregation Block. We first investigate the effectiveness of the LFA block, which is used to capture local information. As shown in Table. VII Row 2, the performance with the MLP-based LFA block is 91.6%/93.1% in terms of mAcc/OA, which is lower than that with the initial LFA block (92.4%/94.0%). This demonstrates that the GCN-based LFA block plays an important role in our algorithm. We also replaced the multi-scale strategy of the LFA block with the single-scale one. As shown in Table. VII Row 3, the classification performance of the multi-scale strategy is superior (91.3%/92.9%). This suggests that the multi-scale features are beneficial to enhancing the expression of local information, thereby improving the performance of our algorithm. Finally, we replaced the Furthest Point Sampling (FPS) method with random sampling, to investigate the performance of 3DGTN with different sampling approaches. As shown in Table. VII Row 4, the classification accuracy drops slightly with random sampling (92.1%/93.5% in terms of mAcc/OA). This is because compared with random sampling, FPS could maintain the geometric characteristics of the target point cloud better. However, as measured by Frame Per Second, Furthest Point Sampling (15) is more time-consuming than random sampling (23). Therefore, developing an efficient and adaptive sampling method for point cloud processing is one of our future works.

GLocal Feature Learning Block. We conducted a detailed ablation study on the GFL block. As shown in Table. VII, when we removed the GFL block, the performance drops significantly, which demonstrates that the GFL block is essential to our algorithm. Secondly, since the GFL block contains two important mechanisms: PPSA and CSA, we also studied the effectiveness of each mechanism. When the CSA was removed, the classification accuracy (mAcc/OA) drops from 92.4%/94.0% to 91.9%/93.7%. Likewise, when the PPSA

was removed, there is a similar drop (from 92.4%/94.0% to 91.13%/93.2%). These results suggest that both self-attention mechanisms are effective in improving classification performance. Additionally, to further verify the effectiveness of the PPSA mechanism, we replaced it with a regular point-wise self-attention mechanism (treating the F_L as the $Value$ matrix). After replacing, we observe a 0.9% and 0.4% drop in mAcc and OA respectively. This confirms the superiority of our PPSA mechanism.

Multi-level GLocal Feature Concatenation. We studied the effectiveness of the multi-level GLocal feature concatenation. As illustrated in Fig. 1, we concatenate the output feature of each level (module) using a residual connection to generate the multi-level GLocal feature. As shown in Table VII Row 9, when the residual connection was removed, we observed a 0.3% and 0.4% drop in mAcc and OA respectively. This suggests that the multi-level GLocal feature contributes significantly to performance improvement.

V. CONCLUSION

In this paper, we have proposed a hierarchical point cloud representation network for classification and segmentation, named 3DGTN. It is an encoder-decoder architecture. The encoder has a series of GLocal modules for effective feature extraction, each of which consists of two cascaded LFA and GFL blocks. In particular, for the GFL block, we adopt the dual-attention Transformer which combines the PPSA and CSA mechanisms. The novel PPSA mechanism is designed to fuse both global features and local neighborhood information of input points, which is able to improve feature learning ability as GLocal features and mitigate local information loss. The decoder is composed of several MLP layers for efficient point cloud reconstruction. It achieves a better trade-off between accuracy and efficiency than a symmetric decoder.

Extensive experiments on the ModelNet40, ScanObjectNN classification datasets [22], [43], ShapeNet part segmentation dataset [44], and MS-LiDAR semantic segmentation dataset [45] demonstrate the superiority of our method in dealing with both synthetic data and real-scene LiDAR data.

Future Work. Our hierarchical network uses Euclidean distance-based downsampling and neighborhood search methods, which are time-consuming and cannot serve the semantic information extracted by the network very well. Since the attention map in the Transformer contains rich feature relationships, we plan to utilize the attention map for semantic-based point cloud sampling and grouping as a future research project. To this end, the “superpoint” strategy could be a potential solution.

REFERENCES

- [1] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Proc. 31st Int. Conf. Neural Inf. Process. Syst.*, 2017, pp. 6000–6010.
- [2] D. Lu, Q. Xie, M. Wei, L. Xu, and J. Li, “Transformers in 3D point clouds: A survey,” *arXiv:2205.07417*, 2022. [Online]. Available: <http://arxiv.org/abs/2205.07417>
- [3] X. Qiang, W. He, S. Chen, Q. Lv, and F. Huang, “Hierarchical point cloud transformer: A unified vegetation semantic segmentation model for multisource point clouds based on deep learning,” *IEEE Trans. Geosci. Remote Sens.*, vol. 61, pp. 1–16, 2023.

- [4] Z. Luo, Z. Zeng, W. Tang, J. Wan, Z. Xie, and Y. Xu, “Dense dual-branch cross attention network for semantic segmentation of large-scale point clouds,” *IEEE Trans. Geosci. Remote Sens.*, vol. 62, pp. 1–16, 2024.
- [5] Y. Song, F. He, Y. Duan, T. Si, and J. Bai, “LSLPCT: An enhanced local semantic learning transformer for 3D point cloud analysis,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–13, 2022.
- [6] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu, “PCT: Point cloud transformer,” *Comput. Vis. Media.*, vol. 7, no. 2, pp. 187–199, Jun, 2021, doi:10.1007/s41095-021-0229-5.
- [7] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “PointNet: Deep learning on point sets for 3D classification and segmentation,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 77–85.
- [8] X.-F. Han, Z.-Y. He, J. Chen, and G.-Q. Xiao, “3CROSSNet: Cross-level cross-scale cross-attention network for point cloud representation,” *IEEE Robotics Autom. Lett.*, vol. 7, no. 2, pp. 3718–3725, 2022.
- [9] 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 *IEEE Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 8500–8509.
- [10] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2021, pp. 9992–10002, doi:10.1109/ICCV48922.2021.00986.
- [11] H. Zhao, L. Jiang, J. Jia, P. H. Torr, and V. Koltun, “Point transformer,” in *Proc. IEEE Int. Conf. Comput. Vis.*, 2021, pp. 16 259–16 268.
- [12] Y. Gao, X. Liu, J. Li, Z. Fang, X. Jiang, and K. M. S. Huq, “LFT-Net: Local feature transformer network for point clouds analysis,” *IEEE Trans. Intell. Transport. Syst.*, 2022, doi:10.1109/TITS.2022.3140355.
- [13] Z. Liu, X. Yang, H. Tang, S. Yang, and S. Han, “FlatFormer: Flattened window attention for efficient point cloud transformer,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 1200–1211.
- [14] D. Lu, Q. Xie, K. Gao, L. Xu, and J. Li, “3DCTN: 3D convolution-transformer network for point cloud classification,” *IEEE Trans. Intell. Transport. Syst.*, 2022, doi:10.1109/TITS.2022.3198836.
- [15] M. Feng, L. Zhang, X. Lin, S. Z. Gilani, and A. Mian, “Point attention network for semantic segmentation of 3D point clouds,” *Pattern Recognit.*, vol. 107, p. 107446, 2020, doi:10.1016/j.patcog.2020.107446.
- [16] L. Hui, H. Yang, M. Cheng, J. Xie, and J. Yang, “Pyramid point cloud transformer for large-scale place recognition,” in *Proc. IEEE Int. Conf. Comput. Vis.*, 2021, pp. 6098–6107.
- [17] C. Zhang, H. Wan, X. Shen, and Z. Wu, “Patchformer: An efficient point transformer with patch attention,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 11 799–11 808.
- [18] J. Sun, C. Qing, J. Tan, and X. Xu, “Superpoint transformer for 3D scene instance segmentation,” in *AAAI Conf. Artif. Intell.*, vol. 37, no. 2, 2023, pp. 2393–2401.
- [19] D. Robert, H. Raguet, and L. Landrieu, “Efficient 3d semantic segmentation with superpoint transformer,” *arXiv:2306.08045*, 2023. [Online]. Available: <http://arxiv.org/abs/2306.08045>
- [20] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, “Dynamic graph CNN for learning on point clouds,” *ACM Trans. Graph.*, vol. 38, no. 5, pp. 1–12, 2019.
- [21] S. Qiu, S. Anwar, and N. Barnes, “Geometric back-projection network for point cloud classification,” *IEEE Trans. Multimedia*, vol. 24, pp. 1943–1955, 2022.
- [22] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, “3D shapenets: A deep representation for volumetric shapes,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2015, pp. 1912–1920.
- [23] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, “PointNet++: Deep hierarchical feature learning on point sets in a metric space,” in *Proc. 31st Int. Conf. Neural Inf. Process. Syst.*, 2017, p. 5105–5114.
- [24] M. Lin and A. Feragen, “diffconv: Analyzing irregular point clouds with an irregular view,” in *Proc. Eur. Conf. Comput. Vis.* Springer, 2022, pp. 380–397.
- [25] A. Muzahid, W. Wan, F. Sohel, L. Wu, and L. Hou, “Curvenet: Curvature-based multitask learning deep networks for 3D object recognition,” *IEEE CAA J. Autom. Sinica*, vol. 8, no. 6, pp. 1177–1187, 2021, doi:10.1109/JAS.2020.1003324.
- [26] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, and B. Chen, “PointCNN: Convolution on X-transformed points,” in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 31, 2018, pp. 820–830.
- [27] C. Kaul, N. Pears, and S. Manandhar, “FatNet: A feature-attentive network for 3D point cloud processing,” in *Proc. Int. Conf. on Pattern Recogn. (ICPR)*, Jan. 2020, pp. 7211–7218, doi:10.1109/ICPR48806.2021.9412731.

- [28] S. Qiu, S. Anwar, and N. Barnes, "Dense-resolution network for point cloud classification and segmentation," in *Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, Jan. 2021, pp. 3813–3822, doi:10.1109/WACV48630.2021.00386.
- [29] X. Ma, C. Qin, H. You, H. Ran, and Y. Fu, "Rethinking network design and local geometry in point cloud: A simple residual MLP framework," in *Proc. Int. Conf. Learn. Represent. (ICLR)*, 2022.
- [30] R. Zhang, L. Wang, Y. Wang, P. Gao, H. Li, and J. Shi, "Starting from non-parametric networks for 3D point cloud analysis," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 5344–5353.
- [31] H. Ran, J. Liu, and C. Wang, "Surface representation for point clouds," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 18942–18952.
- [32] J. Yang, Q. Zhang, B. Ni, L. Li, J. Liu, M. Zhou, and Q. Tian, "Modeling point clouds with self-attention and gumbel subset sampling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 3323–3332, doi:10.1109/CVPR.2019.00344.
- [33] N. Engel, V. Belagiannis, and K. Dietmayer, "Point transformer," *IEEE Access*, vol. 9, pp. 134 826–134 840, 2021, doi:10.1109/ACCESS.2021.3116304.
- [34] X.-F. Han, Y.-J. Kuang, and G.-Q. Xiao, "Point cloud learning with transformer," *arXiv:2104.13636*, 2021. [Online]. Available: <http://arxiv.org/abs/2104.13636>
- [35] Y. Song, F. He, Y. Duan, T. Si, and J. Bai, "Lslpct: An enhanced local semantic learning transformer for 3D point cloud analysis," *IEEE Trans. Geosci. Remote Sens.*, 2022, doi:10.1109/TGRS.2022.3202823.
- [36] K. Mazur and V. Lempitsky, "Cloud transformers: A universal approach to point cloud processing tasks," in *Proc. IEEE Int. Conf. Comput. Vis. (ICCV)*, Oct. 2021, pp. 10695–10704, doi:10.1109/ICCV48922.2021.01054.
- [37] Y. Xu, T. Fan, M. Xu, L. Zeng, and Y. Qiao, "SpiderCNN: Deep learning on point sets with parameterized convolutional filters," in *Proc. Eur. Conf. Comput. Vis.*, vol. 11212, 2018, pp. 90–105.
- [38] A. Goyal, H. Law, B. Liu, A. Newell, and J. Deng, "Revisiting point cloud shape classification with a simple and effective baseline," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2021, pp. 3809–3820.
- [39] S. Cheng, X. Chen, X. He, Z. Liu, and X. Bai, "PRA-Net: Point relation-aware network for 3D point cloud analysis," *IEEE Trans. Image Process.*, vol. 30, pp. 4436–4448, 2021.
- [40] Y. Pang, W. Wang, F. E. Tay, W. Liu, Y. Tian, and L. Yuan, "Masked autoencoders for point cloud self-supervised learning," in *Eur. Conf. Comput. Vis.* Springer, 2022, pp. 604–621.
- [41] A. Berg, M. Oskarsson, and M. O'Connor, "Points to Patches: Enabling the use of self-attention for 3D shape recognition," in *Int. Conf. Pattern Recognit.* IEEE, 2022, pp. 528–534.
- [42] 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 *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2022, pp. 19313–19322.
- [43] M. A. Uy, Q.-H. Pham, B.-S. Hua, T. Nguyen, and S.-K. Yeung, "Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 1588–1597.
- [44] L. Yi, V. G. Kim, D. Ceylan, I.-C. Shen, M. Yan, H. Su, C. Lu, Q. Huang, A. Sheffer, and L. Guibas, "A scalable active framework for region annotation in 3D shape collections," *ACM Trans. Graph.*, vol. 35, no. 6, pp. 1–12, 2016.
- [45] P. Zhao, H. Guan, D. Li, Y. Yu, H. Wang, K. Gao, J. M. Junior, and J. Li, "Airborne multispectral lidar point cloud classification with a feature reasoning-based graph convolution network," *Int J Appl Earth Obs Geoinf.*, vol. 105, p. 102634, 2021.
- [46] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-CAM: Visual explanations from deep networks via gradient-based localization," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2017, pp. 618–626.
- [47] S. Xie, S. Liu, Z. Chen, and Z. Tu, "Attentional shapecontextnet for point cloud recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4606–4615.
- [48] M. Atzmon, H. Maron, and Y. Lipman, "Point convolutional neural networks by extension operators," *arXiv:1803.10091*, 2018. [Online]. Available: <http://arxiv.org/abs/1803.10091>
- [49] 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 *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2530–2539, doi:10.1109/CVPR.2018.00268.
- [50] W. Wang, R. Yu, Q. Huang, and U. Neumann, "SGPN: similarity group proposal network for 3D point cloud instance segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 2569–2578.
- [51] B. Graham, M. Engelcke, and L. Van Der Maaten, "3D semantic segmentation with submanifold sparse convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 9224–9232.
- [52] W. Wu, Z. Qi, and L. Fuxin, "PointConv: Deep convolutional networks on 3D point clouds," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2019, pp. 9621–9630, doi:10.1109/CVPR.2019.00985.
- [53] X. Liu, Z. Han, Y.-S. Liu, and M. Zwicker, "Point2sequence: Learning the shape representation of 3D point clouds with an attention-based sequence to sequence network," in *AAAI Conf. Artif. Intell.*, vol. 33, no. 01, 2019, pp. 8778–8785.
- [54] —, "Point-Voxel CNN for efficient 3D deep learning," in *Adv. Neural Inf. Process. Syst.*, vol. 32, 2019, pp. 963–973.
- [55] Y. Liu, B. Fan, S. Xiang, and C. Pan, "Relation-shape convolutional neural network for point cloud analysis," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 8895–8904.
- [56] H. Thomas, C. R. Qi, J.-E. Deschard, B. Marcotegui, F. Goulette, and L. J. Guibas, "KPConv: Flexible and deformable convolution for point clouds," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 6411–6420.
- [57] J. Mao, X. Wang, and H. Li, "Interpolated convolutional networks for 3D point cloud understanding," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 1578–1587.
- [58] Y. Liu, B. Fan, G. Meng, J. Lu, S. Xiang, and C. Pan, "Densepoint: Learning densely contextual representation for efficient point cloud processing," in *Proc. IEEE Int. Conf. Comput. Vis.*, 2019, pp. 5239–5248.
- [59] M. Xu, R. Ding, H. Zhao, and X. Qi, "PAConv: Position adaptive convolution with dynamic kernel assembling on point clouds," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2021, pp. 3173–3182.
- [60] C. Wu, J. Zheng, J. Pfommer, and J. Beyerer, "Attention-based point cloud edge sampling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2023, pp. 5333–5343.
- [61] L. Wang, Y. Huang, Y. Hou, S. Zhang, and J. Shan, "Graph attention convolution for point cloud semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 10 296–10 305.
- [62] Z. Jing, H. Guan, P. Zhao, D. Li, Y. Yu, Y. Zang, H. Wang, and J. Li, "Multispectral LiDAR point cloud classification using SE-Pointnet++," *Remote Sens.*, vol. 13, no. 13, p. 2516, 2021.
- [63] K. Xiao, J. Qian, T. Li, and Y. Peng, "Multispectral LiDAR point cloud segmentation for land cover leveraging semantic fusion in deep learning network," *Remote Sens.*, vol. 15, no. 1, p. 243, 2022.
- [64] Q. Zhang, Y. Peng, Z. Zhang, and T. Li, "Semantic segmentation of spectral lidar point clouds based on neural architecture search," *IEEE Trans. Geosci. Remote Sens.*, 2023.



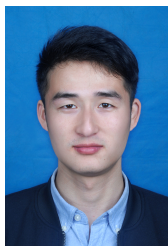
Intelligent Transportation Systems, and ICCV.

Dening Lu received his BSc and MSc degrees in Electrical Engineering, both from the Nanjing University of Aeronautics and Astronautics (NUAA), China in 2018 and 2021, respectively. He is currently pursuing his Ph.D. degree in Systems Design Engineering with the Geospatial Sensing and Data Intelligence Group at the University of Waterloo, Canada. His research interests include 3D point cloud processing and deep learning. He has published papers in IEEE Transactions on Instrumentation and Measurement, IEEE Transactions on



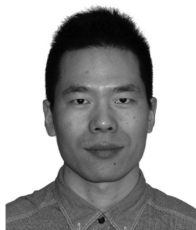
Intelligent Transportation Systems.

Kyle Gao (Graduate Student Member, IEEE) received his Bachelor's degree in mathematics from the University of Waterloo, Canada in 2016 and his Master's degree in physics from the University of Victoria, Canada in 2020, respectively. He is currently pursuing his Ph.D. degree in Systems Design Engineering with the Geospatial Sensing and Data Intelligence Group at the University of Waterloo, Canada. His research interests include computer vision and deep learning. He has published papers in the International Journal of Applied Earth Observation and Geoinformation, Geomatica, and IEEE Transactions on



Qian Xie received his BSc and Ph.D. degrees in Electrical Engineering, both from Nanjing University of Aeronautics and Astronautics (NUAA), China in 2014 and 2021. He is currently a Research Associate in the Department of Computer Science at the University of Oxford, UK. He works in the Cyber-Physical Systems Group under the supervision of Prof. Niki Trigoni and Prof. Andrew Markham. Prior to Oxford, he went to Cardiff University, UK as a joint-trained PhD student in 2019 for 18 months. His research interests are 3D vision, point cloud

processing, deep learning and scene understanding. He publishes extensively in venues and journals such as CVPR, ICCV, ECCV, RA-L, TMM, IJCV, and etc.



Linlin Xu (Member, IEEE) received the B.Eng. and M.Sc. degrees in geomatics engineering from China University of Geosciences, Beijing, China, in 2007 and 2010, respectively, and the Ph.D. degree in geography from the Department of Geography and Environmental Management, University of Waterloo, Waterloo, ON, Canada, in 2014. He is currently a Research Assistant Professor with the Department of Systems Design Engineering, University of Waterloo. He has published various papers on high-impact remote sensing journals and conferences. His

research interests include hyperspectral imaging, synthetic aperture radar and Lidar, machine learning and environmental monitoring.



Jonathan Li (Fellow, IEEE) received the Ph.D. degree in geomatics engineering from the University of Cape Town, South Africa, in 2000. He is currently a Professor of Geomatics and Systems Design Engineering, University of Waterloo, Canada. He has supervised 120+ Masters and PhD students as well as postdoc fellows to completion and co-authored 490+ publications, 320+ of which were published in refereed journals, including IEEE Transactions on Geoscience and Remote Sensing (TGRS), IEEE Transactions on Intelligent Transportation Systems

(TITS), IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing (JSTARS), ISPRS-JPRS, RSE and JAG. He has also published papers in flagship conferences in computer vision and AI, including CVPR, AAAI, and IJCAI. He is the Editor-in-Chief of the International Journal of Applied Earth Observation and Geoinformation (JAG) and the Associate Editor of the TGRS and TITS. His main research interests include AI-based information extraction from earth observation images and LiDAR point clouds as well as 3D vision and GeoAI. He is a Fellow of the Royal Society of Canada (RSC) Academy of Science, the Canadian Academy of Engineering (CAE), and the Engineering Institute of Canada (EIC).