# DLA-Net: Learning Dual Local Attention Features for Semantic Segmentation of Large-Scale Building Facade Point Clouds

Yanfei Su[a], Weiquan Liu[a], Zhimin Yuan[a], Ming Cheng[a,*], Zhihong Zhang[a], Xuelun Shen[a], Cheng Wang[a]

[a]*Fujian Key Laboratory of Sensing and Computing for Smart City, School of Informatics, Xiamen University, Xiamen, China*

## Abstract

Semantic segmentation of building facade is significant in various applications, such as urban building reconstruction and damage assessment. As there is a lack of 3D point clouds datasets related to the fine-grained building facade, we construct the first large-scale building facade point clouds benchmark dataset for semantic segmentation. The existing methods of semantic segmentation cannot fully mine the local neighborhood information of point clouds. Addressing this problem, we propose a learnable attention module that learns Dual Local Attention features, called DLA in this paper. The proposed DLA module consists of two blocks, including the self-attention block and attentive pooling block, which both embed an enhanced position encoding block. The DLA module could be easily embedded into various network architectures for point cloud segmentation, naturally resulting in a new 3D semantic segmentation network with an encoder-decoder architecture, called DLA-Net in this work. Extensive experimental results on our constructed building facade dataset demonstrate that the proposed DLA-Net achieves better performance than the state-of-the-art methods for semantic segmentation.

*Keywords:* semantic segmentation, building facade, self-attention, attentive pooling, DLA-Net

*Corresponding: Ming Cheng
   *Email address:* `chm99@xmu.edu.cn` (Ming Cheng)

## 1. Introduction

Automatic semantic segmentation of building facades is extremely important for urban building modeling and such models have a wide range of applications like urban reconstruction [1],[2],[3],[4] and damage assessment [5]. In the past few years, many works on semantic segmentation of building facades were based on 2D images [6],[7],[8],[9]. However, images acquired by traditional optical imaging-based systems have some intrinsic deficiencies, such as lacking accurate geospatial information, unstable image qualities influenced by illumination conditions, and image distortions caused by camera lens. Furthermore, current semantic segmentation datasets of building facade based on 2D images are insufficient in scale, e.g., the dataset in [9] only contains about 600 images. Some works used Structure from Motion (SfM) to transform the multi-view 2D images into 3D point clouds before building facades semantic segmentation [10],[11]. These methods achieve better segmentation quality and higher speed than those based on 2D images. However, the process of converting 2D images to 3D point clouds by the SfM algorithm is time-consuming, and the generated 3D point clouds are sparse in density and limited with the scale of the area.

In recent years, deep learning has achieved great successes in both image interpretation and 3D point cloud processing. However, one of the main bottlenecks of deep learning is that learning a good model requires sufficient manually labeled training data. To our knowledge, there is a lack of large-scale 3D point clouds datasets of fine-grained building facade, which delays the development of deep learning methods in the tasks of 3D point clouds building facade segmentation.

With the rapid development of 3D sensors, large-scale and highly dense 3D point clouds, which provide rich geometric, shape, and texture information, are easily acquired by Mobile Laser Scanners (MLS) [12],[13] in a short time period. By turning to MLS, we aim at closing this data gap to help unleash
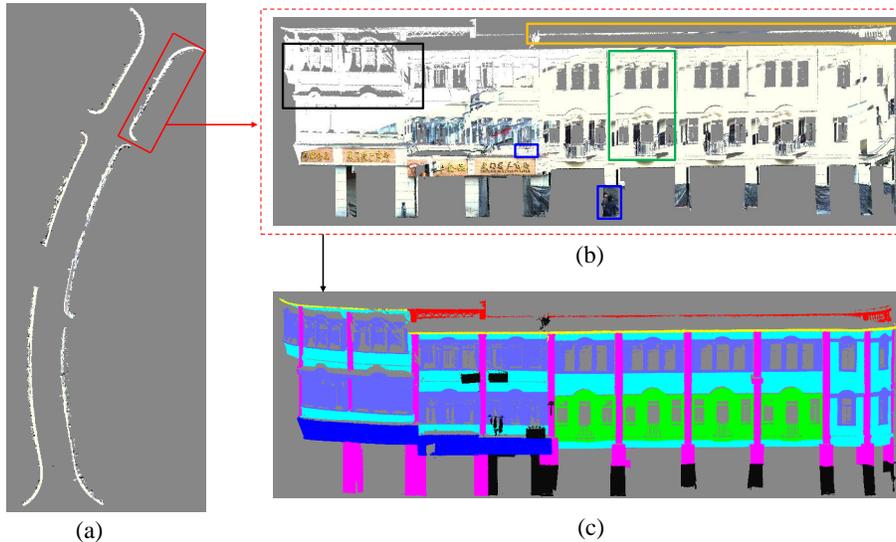
2

**Fig. 1.** Example of semantic labeling of building facade using color mobile LiDAR point clouds: (a) overhead view of the building facades on both sides of the street; (b) enlarged front view of the building facade within the red box in (a); (c) semantic labeling results.

the full potential of deep learning methods for 3D building facade segmentation tasks. In this paper, we construct the first large-scale fine-grained building facade point clouds dataset benchmark obtained by MLS for semantic segmentation, which covers about three kilometers of city street scenes, and over 160 million manually labeled points. A part of the dataset is visualized in Fig. 1. We believe that with such a dataset, the research on building facades semantic segmentation would be promoted. Compared with 2D images, 3D point clouds acquired by MLS can provide real-world coordinates of building facades, and immune to the impact of illumination conditions and image distortions. Although 3D point clouds provide rich scene semantic information, as we can see from Fig. 1(b), segmentation of 3D point clouds is still a challenging task due to the presence of incompleteness (orange bounding box), similar category (green bounding box), color information loss (black bounding box) and occlusion (blue bounding boxes).

In the past few years, a lot of approaches have been proposed for 3D point

3

cloud segmentation, which benefit from the remarkable success of Deep Neural Network (DNN) on irregular point cloud processing. These methods could be roughly divided into three categories [14]: projection-based method [15],[16], [17],[18], discretization-based methods [19],[20],[21],[22], and point-based methods [23],[24],[25], [26]. Both the projection-based and discretization-based methods are computationally expensive to handle large-scale point clouds, which need extra procedures to transform point clouds to a regular representation and project the intermediate segmentation results back to the point clouds. Due to the generation of intermediates, these methods have not fully exploited the underlying geometric and structural information, which inevitably leads to information loss.

Different from the projection-based and discretization-based methods, point-based methods directly working on 3D point clouds have been proposed. Point-Net [23], as the pioneer work, has received much more attention based on which a variety of point-based methods were proposed [27],[28],[29]. Although these methods have achieved promising performances for semantic segmentation, most of them are limited to extremely small 3D point clouds and cannot be directly extended to larger point clouds for two main reasons: (1) These networks usually employ Farthest-Point Sampling (FPS) to obtain keypoint candidates. However, FPS has a time complexity of $O(N^2)$, therefore it is inefficient and time-consuming. (2) Some networks rely on building kernel convolution or graph convolution to learn local features. These methods are also time-consuming and unable to process large-scale point clouds. Lately, RandLA-Net was proposed [26], which utilized random sampling for candidate selection and subsampling. The time complexity of random sampling is $O(1)$ so that it is highly efficient to process large-scale point clouds. In this paper, we also utilize random sampling for candidate selection and subsampling.

Attention-based methods have been flourishing in recent years [30],[31],[32], [33]. The attention mechanism automatically learns important local features by assigning more weight to the key information. Specially, inspired by the success of transformers in natural language processing [34],[35] and image anal-

4

ysis [36],[37], we develop a method for point clouds processing. The core of the transformer is the self-attention operator, which is invariant to permutation and cardinality [37],[38]. So, it is natural to use self-attention as a module in a deep learning network working on 3D point clouds.

In this paper, we aim to design a simple and effective neural network based on an attention module that can directly process large-scale 3D point clouds. The attention module learns Dual Local Attention features from point clouds, which is called DLA module in this paper. The DLA module consists of two blocks, i.e., the self-attention block and the attentive pooling block. First, we aggregate different spatial information to form an enhanced position encoding block, which is embedded in both the self-attention block and the attentive pooling block. Second, the self-attention block concentrates on learning feature representation of the local neighborhoods around each point. Third, the attentive pooling block is used to aggregate local neighboring point feature learned by the self-attention block. The important local features are automatically focused on by the attentive pooling block. Summarily, the proposed DLA module learns more local neighborhood features from point clouds through the self-attention block and the attentive pooling block. Experiments on the building facade dataset show that the network constructed by DLA achieves better segmentation results than the state-of-the-art networks.

In summary, our main contributions in this work include the following:

- We construct the first large-scale fine-grained building facade 3D point clouds dataset benchmark for semantic segmentation. The dataset will be made public soon[1] and we believe that the dataset will further boost the research of deep learning on 3D building facade point clouds.

- We propose an enhanced position encoding block, which aggregates different spatial information to learn more local geometric structure information.

---

[1] https://github.com/suyanfei/DLA-Net

5

- We propose the Dual Local Attention (DLA) module, which consists of two blocks including the self-attention block and the attentive pooling block and is able to capture more local neighborhood features of points. The DLA module can be easily applied to various architectures to explore novel point cloud segmentation networks.

- Extensive experimental results on the building facade dataset demonstrate that the proposed DLA-Net by embedding the DLA module into a standard encoder-decoder architecture achieves state-of-the-art performance.

## 2. Related Work

In this section, we introduce the three learning-based categories of large-scale point clouds segmentation methods, including the projection-based, the discretization-based, and the point-based networks [14].

### 2.1. *Projection-based methods*

Inspired by the success of 2D Convolutional Neural Networks (CNNs), many existing works [16],[39] project 3D point clouds onto 2D images from multiple virtual views to address the task of semantic segmentation. In a related method, Tatarchenko et al. [17] project the local surface geometry around each point onto a tangent plane, forming tangent images that can be processed by 2D convolution. However, the projection step of these multi-view segmentation methods inevitably introduces information loss of the details and as a result, the underlying geometric and structural information is not fully utilized.

### 2.2. *Discretization-based methods*

The discretization-based methods usually transform the point clouds into a discrete representation, such as voxels. In [20], the point clouds are divided into voxels and fed to a fully-3D CNN for voxel-wise segmentation. Some methods use the advantage of 3D CNN for point clouds semantic segmentation [40],[41].

In particular, the Fully-Convolutional Point Network (FCPN) [42] uses 3D convolutions and weighted average pooling to extract features. The discretization-based methods have more flexibility to process the large-scale point clouds. However, the voxelization step inherently introduces discretization artifacts and information loss.

## 2.3. *Point-based methods*

Different from projection-based and discretization-based methods, point-based methods directly work on 3D point clouds. The pioneering work Point-Net [23] learns a spatial encoding for each point using pointwise MLPs and then aggregates all individual point features as a global representation using symmetrical pooling functions. Based on PointNet, a variety of point-based methods were proposed including [24],[27],[26], etc. Overall, these methods can be roughly divided into pointwise MLP methods, point convolution methods, RNN-based methods, graph-based methods and attention-based methods.

**Pointwise MLP methods**. PointNet++ [24], a hierarchical spatial structure, learns a feature for each point by aggregating the information from local neighboring points. Hu et al. [26] propose an efficient and lightweight network called RandLA-Net for large-scale point cloud segmentation. RandLA-Net utilizes the random point sample method to achieve remarkably high efficiency in terms of memory and computation. Meanwhile, a local feature aggregation module is further proposed to capture and preserve geometric features.

**Point convolution methods**. A handful of approaches are based on effective convolution operators for point clouds. PointCNN [27] transforms neighboring points to the canonical order, which enables traditional convolution to play a normal role. KPConv [29] proposes a spatially deformable point convolution with any number of kernel points which alleviates both varying densities and computational cost, outperforming all associate methods on point clouds segmentation tasks.

**Graph-based and RNN-based methods**. To capture the underlying shapes and local structure features from point clouds, graph networks and Re-

7

current Neural Networks (RNN) networks have also been used for semantic segmentation of point clouds. DGCNN [43] proposes an EdgeConv module, which generates edge features that describe the relationships between a point and its neighbors. RSNet [44] proposes a lightweight local dependency module and utilizes a slice pooling layer to project the feature of unordered points onto an ordered sequence of feature vectors.

**Attention-based methods**. In recent years, the powerful attention mechanism has attracted more and more attention [32],[33]. Wang et al. [30] first use the graph pointnet module based on graph attention block to dynamically compose and update each point representation within the local point cloud structure, then resorts to the spatial-wise and channel-wise attention strategies to exploit the point cloud global structure. In particular, transformers and self-attention have shown remarkable performance in machine translation and natural language processing [34],[35]. Inspired by transformer, many works applied a self-attention network into 2D image recognition [37],[36]. So, in this paper, we try to apply the self-attention network to 3D point clouds of building facades for semantic segmentation.

## 3. Method

In this section, we begin with a brief overview of transformer and self-attention. Then we detail the position encoding block, the self-attention block, and the attentive pooling block for 3D point clouds processing. Lastly, we present our network architecture for 3D semantic segmentation.

### 3.1. *Background*

The deep learning networks based on transformer and self-attention have become a popular and efficient method in the field of natural language processing and 2D image analysis [34],[35],[36],[37],[38]. An attention mechanism uses input-dependent weights to linearly combine the inputs. In [38], the pairwise self-attention is explored. Mathematically, let $\chi = \{x_i\}$ be a set of feature

vectors, the pairwise self-attention has the following form:

$$y_i = \sum_{x_j \in \chi} softmax\left( \eta \Big( \delta \big( \alpha(x_i), \beta(x_j) \big) + \rho \Big) \right) \odot \gamma(x_j), \tag{1}$$

where $y_i$ is the output feature corresponding to $x_i$, $\odot$ is the Hadamard product, $i$ is the spatial index of feature vector $x_i$. $\alpha$, $\beta$, and $\gamma$ are trainable transformations such as linear projections or MLPs. $\delta$ is a relation function including summation, subtraction, concatenation, Hadamard product, or dot product. The experimental results in [38] show that summation, subtraction, and Hadamard product have the same performance and are better than concatenation and dot product. $\rho$ is a position encoding function and $\eta$ is a mapping function such as MLPs.

### 3.2. *Position encoding block*

Position encoding plays a critical role in networks that are based on transformer and self-attention. In natural language processing, the sine and cosine functions are usually used as position encoding to provide effective position information for sequences [34]. In 2D image processing, the relative position of 2D coordinates is used as position encoding to augment image features [38]. In 3D point clouds processing, the 3D point coordinates themselves are a natural candidate for position encoding. We design an enhanced position encoding block based on relative point position, which can better explore the local geometric structure. Given the input clouds with $N$ points, whose $xyz$ positions are denoted as $P = \{p_1, \ldots, p_i, \ldots, p_N\} \subset \mathbb{R}^3$. For a center point $p_i$, its neighboring points, denoted as $\{p_i^1, \ldots, p_i^k, \ldots, p_i^K\} \subseteq P$, are usually gathered by the K-nearest neighbors (KNN) algorithm, which is based on the point-wise Euclidean distances. Our position encoding can be expressed as follows:

$$c_i^k = MLPs\Big( (p_i - p_i^k) \oplus \| p_i - p_i^k \| \Big), \tag{2}$$

where $\oplus$ is the concatenation operation, $\| \cdot \|$ calculates the Euclidean distance between the neighboring and the center points, and MLPs consist of two linear
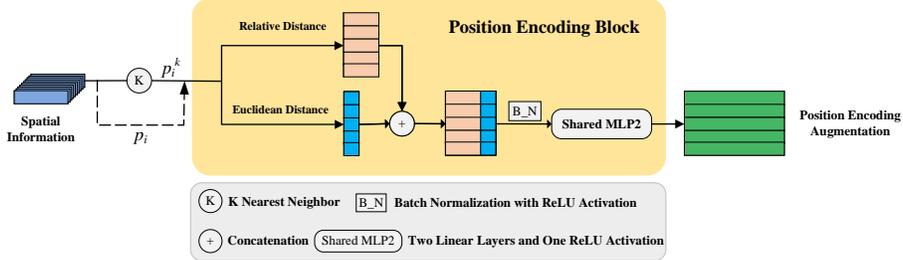
**Fig. 2.** Structure of the position encoding block.

transformations with a ReLU activation. The visual interpretation of the position encoding block is shown in Fig. 2. Through the position encoding block, the spatial information of point clouds will be enhanced, so that helps the network to learn local features and obtain good performance in practice. Note that the aggregation features need to be processed by batch normalization with ReLU activation operation.

### 3.3. Self-attention block

Self-attention is a natural choice for processing point clouds, which are essentially sets of irregularly embedded metric space. As mentioned in Section 3.1, our self-attention block is based on pairwise self-attention. We choose subtraction as the relation function and add a position encoding to both the mapping function $\eta$ and the transformed features $\gamma$. Our self-attention block is illustrated in Fig. 3. Let the feature of the center point $p_i$ be $f_i$, and the feature of its neighboring points $p_i^k$ be $f_i^k$. $f_i, f_i^k \in \mathbb{R}^d$ ($d$ is feature channels). Our self-attention block is represented as follows:

$$F_i = \sum_{k=1}^{K} softmax\Big(\eta\big(\alpha(f_i) - \beta(f_i^k) + c_i^k\big)\Big) \odot \big(\gamma(f_i^k) + c_i^k\big), \tag{3}$$

where $F_i$ is the output feature, $\odot$ is the Hadamard product. $\alpha$, $\beta$, and $\gamma$ are MLP with one linear layer, respectively. The mapping function $\eta$ is an MLP including two linear transformations with a ReLU activation. $c_i^k$ is the position encoding , which comes from Eq. 2. The output of the self-attention block $F_i$ is the new set of neighboring features, which explicitly encodes the local

10

**Fig. 3.** Structure of the self-attention block.

geometric structures for the center point $p_i$. The output feature $F_i$ also needs to be processed by batch normalization with ReLU activation operation.

### 3.4. *Attentive pooling block*

It has been proved in recent work [45] that the attention mechanism automatically learns local features. We continually turn to the attention mechanism to further explore the local features. Our structure of the attentive pooling block is shown in Fig. 4.

Given the output of the self-attention block $F_i$ and position encoding $c_i^k$, we use a relation function to aggregate them. It is formally defined as follows:

$$\hat{F}_i = \eta(F_i, c_i^k). \tag{4}$$

In the attentive pooling block, we use the concatenation $\oplus$ as the relation function $\eta$.

To get an attention weight for the aggregation feature $\hat{F}_i = \{\hat{f}_i^1, \ldots, \hat{f}_i^k, \ldots, \hat{f}_i^K\}$, a shared function $\varphi(.,.)$ is designed, which consists of a shared MLP followed by softmax. The learnable weight $W_i$ in the shared MLP assigns a unique attention score to each feature. The learned attention scores can be regarded as a soft
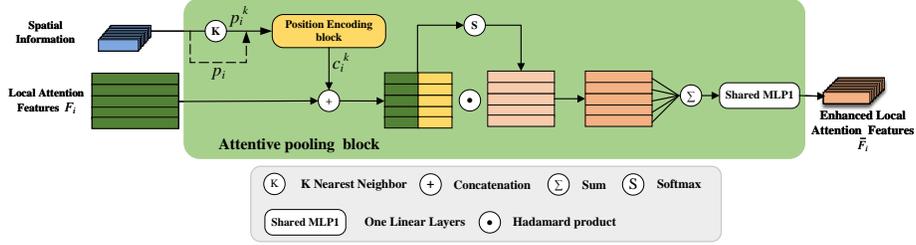
11

**Fig. 4.** Structure of the attentive pooling block.

mask that automatically selects the important features. Then these features are weighted and summed to generate a new informative feature vector $\bar{F}_i$. This process is represented as follows:

$$\bar{F}_i = \sum_{k=1}^{K} \varphi(\hat{f}_i^k, W_i) \odot \hat{f}_i^k, \qquad (5)$$

where $\odot$ is the Hadamard product.

### 3.5. Architecture of DLA-net

The architecture of the proposed DLA module is shown in Fig. 5. The inputs of DLA include the spatial information and the features learned previously. The spatial information is utilized to construct the point encoding block which is embedded in the self-attention block and the attentive pooling block. The features learned previously are used by the DLA module for exploring the local information. Inspired by the successful ResNet [46], we stack self-attention block and attentive pooling block with a skip connection as a residual module. The DLA residual module is simple but effective and achieves state-of-the-art performance on our constructed large-scale building facade dataset.

We embed the proposed DLA residual module into a standard encoder-decoder architecture, resulting in the new segmentation network, DLA-Net. The complete architecture of DLA-net is shown in Fig. 6. The input of the DLA-Net is a large-scale point cloud of size $N \times d_{in}$ where $N$ is the number of points and $d_{in}$ is the input feature dimension. Each point is presented by its 3D coordinates and color information, i.e., $d_{in} = 6$. First, each point is extracted by a fully
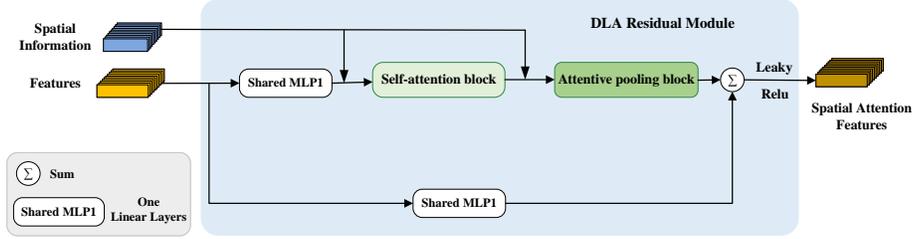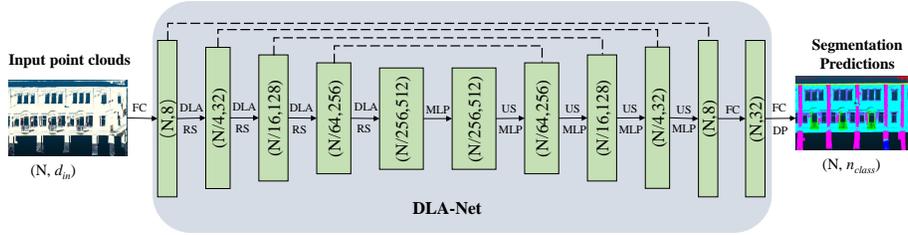
12

**Fig. 5.** Architecture of the DLA.



**Fig. 6.** Architecture of the DLA-Net, (N,$d_{in}$) represents the number of points and feature dimension respectively of the input point clouds; FC: Fully Connected layer; DLA: Dual Local Attention residual module; RS: Random Sampling, MLP: shared Multi-Layer Perceptron; US: Up-sampling; (N,$n_{class}$) represents the number of points and classes respectively of the output point clouds.

connected layer, the dimension is unified to 8 as input to the network. The network consists of four encoder layers and four decoder layers. The encoder layers are used to progressive encode the features in which the DLA module and random sampling operation are embedded. The number of points is gradually decreased from $N$ to $\dfrac{N}{256}$, while the feature dimension is increased from 8 to 512. Second, four decoder layers are used to decode the feature. The decoded features are upsampled through a nearest-neighbor interpolation, and further concatenated with the intermediate feature map produced by encoder layers through skip connections. At last, two shared fully connected layers and a drop layer with ratio 0.5 are used to predict the semantic labels. The outputs of DLA-Net are the predicted semantic labels of all points, with a size of $N \times n_{class}$, where $n_{class}$ is the number of classes.

13

## 4. Experiments

### 4.1. Implementation details

Our experiments are implemented in Tensorflow on a server with an Intel (R) Xeon (R) E5-2683 CPU, 64GB of RAM, an NVIDIA Titan X GPU, CUDA10.0, and cuDNN v7. We use the Adam optimizer with an initial learning rate of $10^{-2}$. The DLA-Net is trained for 100 epochs, with the learning rate dropped by 5% after each epoch. The cross-entropy loss is used for training. The batch size is set as 6. The number of nearest points $K$ is set as 16. A fixed number of points (40960) are sampled from each point cloud and fed into the network while training. The whole raw point cloud is input during testing.

### 4.2. Dataset

In this paper, the fine-grained building facade point clouds dataset was acquired along the National Road and Siming South Road in Xiamen, China, by a RIEGL VMX-450 mobile LiDAR system. The RIEGL VMX-450 system [47] smoothly integrates two RIEGL VQ-450 laser scanners, a global navigation satellite system (GNSS) antenna, an inertial measurement unit (IMU), a distance measurement indicator (DMI), and four high-resolution digital cameras (see Fig. 7). This integrated set was mounted on the roof of a minivan with an average speed of 40-50 km/h. After data acquisition, RiProcess, a post-process software released by RIEGL corporation, is used to calibrate the images with point clouds for the generation of colorized building facade point clouds. The scanning frequency of the VMX-450 scanning system is up to 400Hz, the pulse frequency is 1100kHz, the maximum scanning distance is up to 800m, the scanning accuracy is up to 5 mm, and the positioning accuracy is up to 5 cm. Therefore, the data obtained by the RIEGL VMX-450 onboard mobile laser scanning system is high-resolution 3D laser point clouds data.

The colorized building facade point clouds dataset has a point density of about 7,000 points/$m^2$ and covers two roads sections of total about 3,000 meters. In a real city street scene, the point clouds acquired by the RIEGL VMX-450 system have many non-building facade points such as ground, pedestrians,

14

**Fig. 7.** Illustration of RIEGL VMX-450 mobile LiDAR system and its configurations.

vehicles, trees, outlier, and so forth. So it is necessary to preprocess the data by manually removing non-building facade point clouds. The existence of these points inevitably interferes with the collection of building facade points and causes incomplete data. After data preprocessing, we divide the dataset into 6 areas according to the style and location of the building facade. We use the CloudCompare tool to visualize a representative building facade for each area (see Fig. 8). As seen from Fig. 8, the building facade in the dataset is the typical commercial street scene. Tab. 1 details the number of points in each area and the total of six areas.

**Tab. 1.** The number of points in each area and the total of six areas.

|        | Area 1 | Area 2 | Area 3 | Area 4 | Area 5 | Area 6 |
|--------|--------|--------|--------|--------|--------|--------|
| Points | 33,206,410 | 30,063,042 | 14,790,151 | 24,954,180 | 22,785,852 | 32,959,555 |
| Total  | **158,759,190** | | | | | |

**Tab. 2.** The number of points in each category.

|        | Balustrade | Balcony | Advboard | Wall | Eave | Column | Window | Clutter |
|--------|-----------|---------|----------|------|------|--------|--------|---------|
| Points | 1,535,374 | 7,169,298 | 22,190,094 | 51,651,937 | 5,744,710 | 40,034,103 | 26,287,509 | 4,146,165 |
| Total  | **158,759,190** | | | | | | | |

(a) Area 1

(b) Area 2

(c) Area 3

(d) Area 4

(e) Area 5

(f) Area 6

**Fig. 8.** Illustration of a representative building facade for each area.
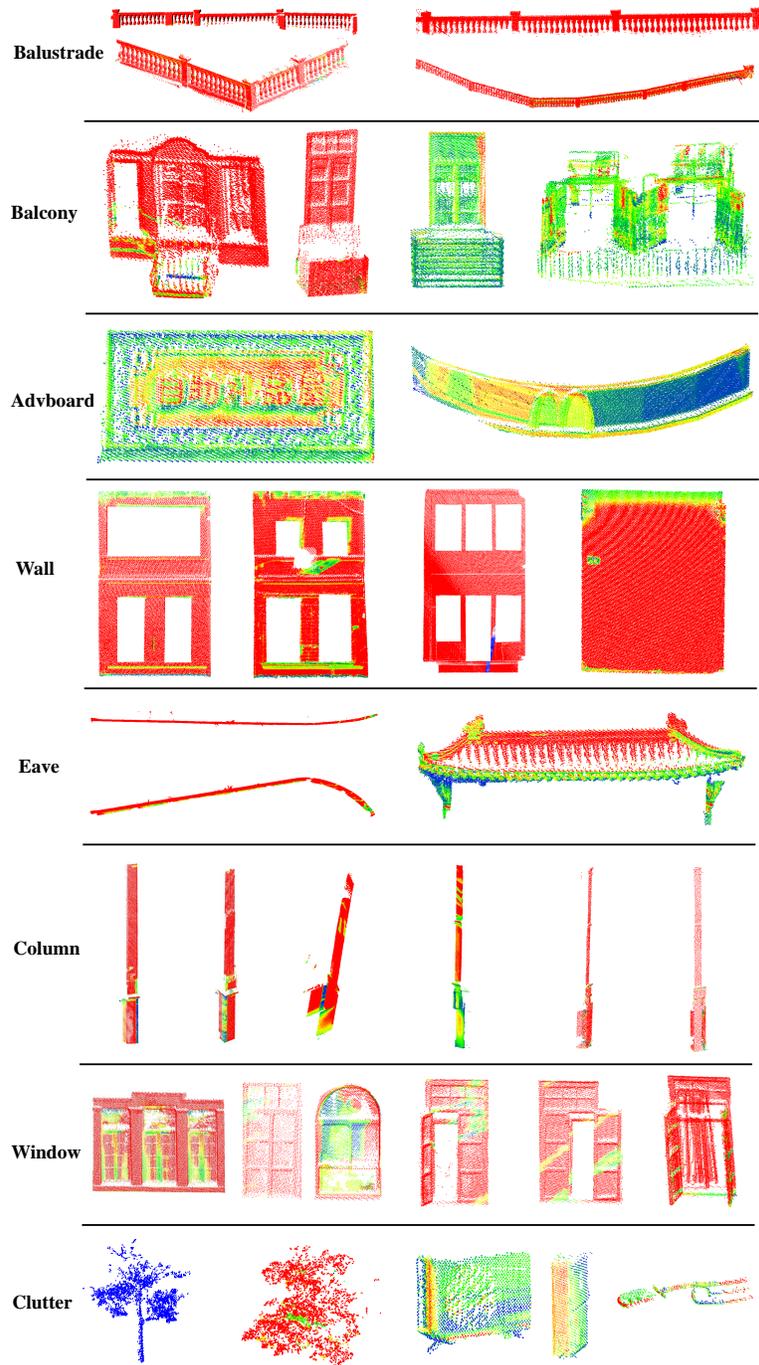
16

**Fig. 9.** Illustration of each category from different perspectives and with different forms.

we construct the ground truth for the dataset by manually and thoroughly classifying all points into the following eight categories: balustrade, balcony, advboard (i.e., advertising board), wall, eave, column, window, and clutter. Tab. 2 details the number of points in each category. The statistical data show the challenge of category imbalance. In order to better intuitively understand each category, we selectively visualize each category from different perspectives and with different forms in Fig. 9. Specially, we regard convex parts and windows as balconies when they are vertically distributed in space. We regard the air conditioner external unit hanging on the wall, weeds and bushes growing on the eaves, wire boxes and street lamps hanging on the column as clutter. The samples of each category are diverse in shape, which shows the richness of our building facade dataset.

Challenges, such as object occlusions, incompleteness, category imbalance, and diversity of category samples, commonly exist in the building facade dataset. These challenges make the semantic segmentation of building facade point clouds a difficult task.

### 4.3. Evaluation on the dataset

To fully evaluate our DLA-Net on building facade dataset for semantic segmentation, we used two modes: (a) Area 1 is used for test and Areas 2 to 6 are used for train, similarly, each area is used for test in turn and others are used for train. (b) k-fold cross-validation (k=6). We use overall pointwise accuracy (OA), mean classwise intersection over union (mIoU), and mean of classwise accuracy (mAcc) as the standard metrics. In addition, 8 methods including PointNet [23], PointNet++ [24], DGCNN [43], ELGS [30], RSNet [44], PointCNN [27], KPConv [29], RandLA-Net [26] are used for the referred approaches which are retrained on our dataset.

The quantitative results of all the referred methods tested on each area are reported in Tab. 3 and Tab. 4. As seen from Table. 3, when Area 1 is used for the test, our method performs better than others on mAcc, and only 0.1% less than KPConv on mIoU. Our DLA-Net achieves better performance than

**Tab. 3.** Quantitative results of different approach on the building facade dataset (Area 1, Area 2, and Area 3 take turns as the test ).

| methods | Area 1–test | | | Area 2–test | | | Area 3–test | | |
|---|---|---|---|---|---|---|---|---|---|
| | OA(%) | mIoU(%) | mAcc(%) | OA(%) | mIoU(%) | mAcc(%) | OA(%) | mIoU(%) | mAcc(%) |
| PointNet [23] | 55.5 | 35.0 | 51.7 | 62.3 | 41.8 | 57.6 | 58.7 | 36.2 | 54.3 |
| PointNet++ [24] | 58.6 | 38.5 | 56.3 | 60.7 | 34.4 | 53.3 | 62.9 | 39.4 | 57.8 |
| DGCNN [43] | 63.1 | 42.1 | 57.7 | 69.4 | 51.0 | 66.7 | 66.5 | 45.2 | 63.1 |
| ELGS [30] | 69.1 | 51.5 | 62.1 | 68.9 | 47.0 | 58.3 | 65.5 | 39.3 | 55.6 |
| RSNet [44] | 68.7 | 50.0 | 65.4 | 69.0 | 50.2 | 64.3 | 70.0 | 47.6 | 65.3 |
| PointCNN [27] | 77.0 | 55.9 | 68.8 | 79.6 | 52.0 | 65.7 | 79.1 | 53.4 | 68.3 |
| KPConv [29] | **85.7** | **68.3** | 78.3 | **84.7** | 66.8 | 77.8 | **83.5** | 56.9 | 67.1 |
| RandLA-net [26] | 80.0 | 63.0 | 77.4 | 81.6 | 62.0 | 77.2 | 80.0 | 55.8 | **73.1** |
| **DLA-Net(ours)** | 83.9 | 68.2 | **81.4** | 84.3 | **66.9** | **81.5** | 83.1 | **58.4** | 71.9 |

**Tab. 4.** Quantitative results of different approach on the building facade dataset (Area 4, Area 5, and Area 6 take turns as the test ).

| methods | Area 4–test | | | Area 5–test | | | Area 6–test | | |
|---|---|---|---|---|---|---|---|---|---|
| | OA(%) | mIoU(%) | mAcc(%) | OA(%) | mIoU(%) | mAcc(%) | OA(%) | mIoU(%) | mAcc(%) |
| PointNet [23] | 51.3 | 27.3 | 39.4 | 54.3 | 30.8 | 45.5 | 47.3 | 24.7 | 36.2 |
| PointNet++ [24] | 44.2 | 21.8 | 33.4 | 63.4 | 39.4 | 54.5 | 57.2 | 32.3 | 45.0 |
| DGCNN [43] | 60.0 | 35.1 | 49.4 | 66.4 | 41.0 | 55.5 | 55.4 | 29.7 | 42.3 |
| ELGS [30] | 52.8 | 28.8 | 45.2 | 69.8 | 51.6 | 63.2 | 62.3 | 36.6 | 49.2 |
| RSNet [44] | 56.5 | 34.5 | 50.4 | 69.1 | 43.5 | 58.0 | 60.5 | 36.2 | 52.8 |
| PointCNN [27] | 70.1 | 41.3 | 53.2 | 78.3 | 50.7 | 64.4 | 65.6 | 38.4 | 52.1 |
| KPConv [29] | 77.1 | 51.2 | 62.7 | 77.9 | 52.4 | 66.8 | **78.4** | **51.7** | **63.0** |
| RandLA-net [26] | **79.4** | 53.9 | **69.4** | 79.8 | 60.3 | 71.0 | 72.1 | 47.7 | 62.2 |
| **DLA-Net(ours)** | 78.5 | **54.3** | 69.1 | **81.2** | **63.6** | **74.0** | 76.0 | 50.4 | 62.2 |

state-of-the-art methods on mIoU when the test set are Area 2 and Area 3. As seen from Tab. 4, when the test set is Area 4, our method has the best mIoU among all the methods. DLA-Net outperforms all prior models according to all metrics when Area 5 is used for the test. In general, DLA-Net achieves better performance on at least one metric except for Area 6.

**Tab. 5.** Semantic segmentation results on the our data set, evaluated with 6-fold cross-validation.

| methods | OA(%) | mIoU(%) | mAcc(%) | Abalustrade | balcony | advboard | wall | eave | column | window | clutter |
|---|---|---|---|---|---|---|---|---|---|---|---|
| PointNet [23] | 53.3 | 32.8 | 46.9 | 12.2 | 15.6 | 45.3 | 31.8 | 39.4 | 56.1 | 30.7 | 31.4 |
| PointNet++ [24] | 57.8 | 34.3 | 50.0 | 39.4 | 12.5 | 41.5 | 46.0 | 45.9 | 44.8 | 36.7 | 8.5 |
| DGCNN [43] | 61.8 | 40.2 | 54.4 | 14.8 | 22.9 | 51.0 | 42.2 | 55.0 | 60.6 | 40.5 | **34.3** |
| ELGS [30] | 63.0 | 40.6 | 54.8 | 18.2 | 21.3 | 56.4 | 46.4 | 51.4 | 58.0 | 43.4 | 29.2 |
| RSNet [44] | 63.9 | 42.1 | 56.0 | 19.7 | 21.8 | 62.7 | 44.4 | 55.6 | 63.2 | 43.4 | 25.7 |
| PointCNN [27] | 75.0 | 48.6 | 62.1 | 43.6 | 9.9 | 70.5 | 60.3 | 58.6 | 76.9 | 53.4 | 15.6 |
| KPConv [29] | **81.3** | 58.0 | 68.0 | 49.7 | 35.7 | 77.5 | **67.9** | **58.7** | **83.7** | 60.4 | 30.4 |
| RandLA-net [26] | 78.6 | 56.0 | 68.6 | 51.4 | 36.4 | 78.3 | 62.9 | 58.5 | 76.9 | 58.8 | 24.9 |
| **DLA-Net(ours)** | 81.0 | **59.7** | **71.6** | **58.2** | **43.3** | **79.1** | 67.5 | 58.0 | 79.0 | **62.0** | 30.2 |

Then, we report the quantitative results of all the referred methods evaluated with 6-fold cross-validation in Tab. 5. As seen from Tab. 5, DLA-Net has the best mIoU and mAcc among all these methods and outperforms the prior state-of-the-art by nearly 3% in mIoU and 4.4% in mAcc. As for OA, it is slightly lower than KPConv, but it is better than all the other compared methods. DLA-Net also achieves the best performance on 4 categories, including balustrade, balcony, advboard, and window.

In the Section 3.2 and 3.3. The batch normalization with ReLU activation plays an important part in the position encoding block and the self-attention block. We define $\times$ as without batch normalization and ReLU activation, and $\sqrt{}$ as batch normalization with ReLU activation. The experiments are conducted on the building facade dataset, tested on Area 1.

As we can see from Tab. 6, when batch normalization with ReLU activation operation is no used in both the Position Encoding Block and self-attention Block, DLA-Net achieves poor performance. When batch normalization with ReLU activation operation is used in the Position Encoding Block

**Tab. 6.** The results of DLA-net by whether to use batch normalization with ReLU activation in the Position Encoding Block and self-attention block.

| Position encoding block | Self-attention block | OA(%) | mIoU(%) | mAcc(%) |
|:---:|:---:|:---:|:---:|:---:|
| × | × | 82.2 | 64.0 | 76.0 |
| × | √ | 83.4 | 66.2 | 78.0 |
| √ | × | 81.5 | 64.8 | 78.5 |
| √ | √ | **83.9** | **68.2** | **81.4** |

or self-attention Block, DLA-Net performs better than not used. Only When batch normalization with ReLU activation operation is used in both the Position Encoding Block and self-attention Block, DLA-Net outperforms all prior models.

The presence of color information is helpful to improve the accuracy of semantic segmentation. Tab. 7 presents the quantitative results of DLA-Net with respect to the different types of input point clouds. The experiments are conducted on the building facade dataset, evaluated with 6-fold cross-validation. When DLA-net is trained given both point coordinates and RGB information, the network achieves significantly better segmentation accuracy.

**Tab. 7.** Quantitative results of DLA-Net on our building facade dateset.

| methods | OA(%) | mIoU(%) | mAcc(%) | balustrade | balcony | advboard | wall | eave | column | window | clutter |
|:---|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| DLA-Net(w/o RGB) | 80.0 | 57.6 | 69.4 | 55.1 | 39.3 | 76.8 | 66.2 | 56.4 | 78.9 | 60.4 | 27.4 |
| DLA-Net(w RGB) | **81.0** | **59.7** | **71.6** | **58.2** | **43.3** | **79.1** | **67.5** | **58.0** | **79.0** | **62.0** | **30.2** |

Lastly, the segmentation results visualization examples of several typical buildings facade are shown in Fig. 10. The figure qualitatively compares the semantic segmentation obtained by PointNet [23], ELGS [30], KPConv [29], RandLA-Net [26], and our DLA-Net. We can see that objects such as the clutter hung on the wall, balcony and window embedded in the wall, clutter cover on the column are quite difficult to segment. As can be seen from Fig. 10, it is obvious that the segmentation results of PointNet and ELGS are the worst.
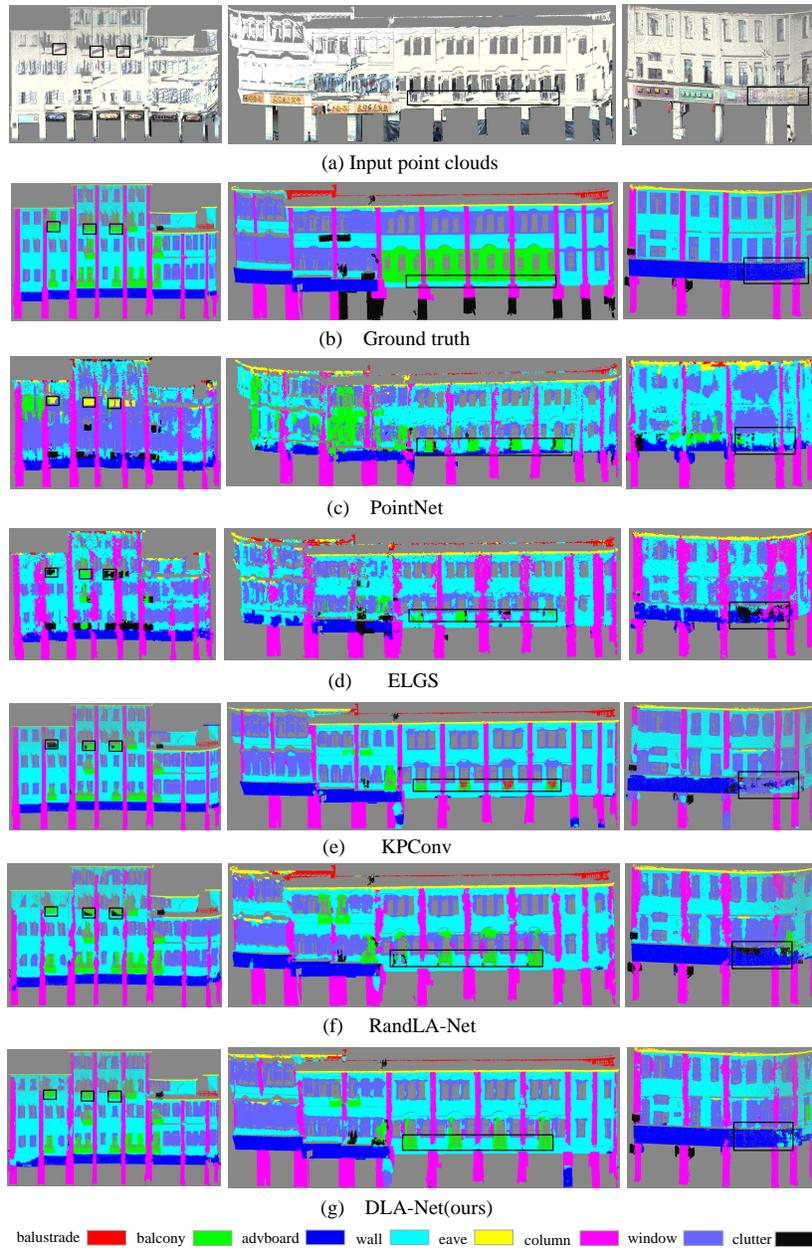
(a) Input point clouds

(b)   Ground truth

(c)   PointNet

(d)   ELGS

(e)   KPConv

(f)   RandLA-Net

(g)   DLA-Net(ours)

balustrade ■  balcony ■  advboard ■  wall ■  eave ■  column ■  window ■  clutter ■

**Fig. 10.** Qualitative comparison on the three scenes of building facade dataset. Different colors denote different categories. These scenes contain 8 categories. Black boxes highlight some examples where our method performs better than others.

We use black bounding boxes to highlight some examples where our method performed significantly better than the competitors. In the first building facade, the convex part of the balcony is very similar to the air conditioner external unit in the spatial shape but different in size making the semantic segmentation a real challenge. So, PointNet [23], ELGS [30], KPConv [29], and RandLA-Net [26] mistake the convex part of the balcony as clutter. The proposed DLA-Net can identify the convex part of the balcony more accurately than these networks. In the second building facade, all the networks failed to separate the clutter from the column, because the clutter is covered on the column. PointNet [23], ELGS [30], KPConv [29], and RandLA-Net [26] also did not fully identify the convex part of the balcony. This is because the convex part of the balcony is diverse in shape, which is also similar to the balustrade. However, the proposed DLA-Net can still accurately identify the convex part of the balcony. In the third building facade, the output of DLA-Net on windows is more regular than that of PointNet [23], ELGS [30], KPConv [29], and RandLA-Net [26], and DLA-Net performs better on advboard than others.

### 4.4. Ablation study

The effectiveness of DLA-Net is verified by the experimental results on the building facade dataset. To better understand the DLA-Net, we further evaluate it and conduct the following three groups of experiments. we also use OA, mIoU and mAcc as the metrics.

#### 4.4.1. Ablation study on position encoding block

The following ablation studies are conducted to study the impacts of the position encoding block in our framework. Position encoding block is a combination of different spatial information. So, we explore the effects of position encoding composed of different spatial information in our DLA-Net. The experiments are conducted on the building facade dataset, evaluated with 6-fold cross-validation. There are six combinations as follows:

(1) Position encoding block consists of the neighboring points $p_i^k$ only.

(2) Position encoding block consists of the relative position $p_i - p_i^k$ only.

(3) Position encoding block consists of the relative position $p_i - p_i^k$ and Euclidean distance $\| p_i - p_i^k \|$.

(4) Position encoding block consists of the points $p_i$, the relative position $p_i - p_i^k$, and Euclidean distance $\| p_i - p_i^k \|$.

(5) Position encoding block consists of the neighboring points $p_i^k$, the relative position $p_i - p_i^k$, and Euclidean distance $\| p_i - p_i^k \|$.

(6) Position encoding block consists of the points $p_i$, the neighboring points $p_i^k$, the relative position $p_i - p_i^k$, and Euclidean distance $\| p_i - p_i^k \|$.

**Tab. 8.** The results of DLA-net by concatenating the different spatial information as the position encoding block.

| Position encoding block | OA(%) | mIoU(%) | mAcc(%) |
|---|---|---|---|
| $p_i^k$ | 69.5 | 47.0 | 61.9 |
| $p_i - p_i^k$ | 80.8 | 59.3 | 71.1 |
| $(p_i - p_i^k) \oplus \| p_i - p_i^k \|$ (**ours**) | **81.0** | **59.7** | **71.6** |
| $p_i \oplus (p_i - p_i^k) \oplus \| p_i - p_i^k \|$ | 80.9 | 58.9 | 70.0 |
| $p_i^k \oplus (p_i - p_i^k) \oplus \| p_i - p_i^k \|$ | 80.6 | 59.2 | 71.2 |
| $p_i \oplus p_i^k \oplus (p_i - p_i^k) \oplus \| p_i - p_i^k \|$ | 80.4 | 58.5 | 69.9 |

The results are shown in Tab. 8. We can see that: 1) Only using the neighboring points $p_i^k$ as the position encoding block is not going to work very well. 2) The combinations based on the relative point position used for the position encoding block get good performance. As we can see from this, the relative position plays an important role in the position encoding block. 3) When the relative position and Euclidean distance are concatenated as the position encoding block, our DLA-Net achieves the best performance. 4) Base on the relative position and Euclidean distance, if the points $p_i$, or the neighboring points $p_i^k$, or together, is concatenated as the position encoding block, the segmentation result is worse than 3). It is because too much spatial information interferes with each other.

### 4.4.2. *Ablation study on self-attention block*

In section 3.3, the position encoding block is added to the mapping function and the transformed features. In this section, we investigate the location of the position encoding block in the self-attention block. The experiments are conducted on the building facade dataset, tested on Area 1.

**Tab. 9.** Ablation study on self-attention pooling block.

| Self-attention block | OA(%) | mIoU(%) | mAcc(%) |
|---|---|---|---|
| None | 83.3 | 66.1 | 77.5 |
| Add only to the mapping function block | 83.5 | 67.6 | 80.5 |
| Add only to transformed features | 83.3 | 66.4 | 79.2 |
| **Ours** | **83.9** | **68.2** | **81.4** |

The results are shown in Tab. 9. It can be viewed that without position encoding block in the self-attention block, the performance of DLA-Net drops. When the position encoding block is added only to the mapping function or only to the transformed features, the performance of DLA-Net is also not good. Only when the position encoding block is both added to the mapping function and the transformed features, the self-attention block perform well.

### 4.4.3. *Ablation study on attentive pooling block*

We verify the effectiveness of attentive pooling block from many aspects, including removing all of it, only removing the position encoding block of it, replace it with max-pooling or avg-pooling. The experiments are conducted on the building facade dataset, tested on Area 1.

As we can see from Tab. 10, the results prove two points: 1) the attentive pooling block plays an important role in DLA; 2) the position encoding block plays an important role in the attentive pooling block.

**Tab. 10.** Ablation study on attentive pooling block.

| Attentive pooling block | OA(%) | mIoU(%) | mAcc(%) |
|---|---|---|---|
| Remove all | 83.1 | 64.7 | 76.2 |
| Remove position encoding block | 82.9 | 65.1 | 77.2 |
| Replace with max-pooling | 79.2 | 62.1 | 77.9 |
| Replace with avg-pooling | 77.8 | 61.0 | 76.7 |
| **Ours** | **83.9** | **68.2** | **81.4** |

## 5. Conclusion

In this paper, we construct the first large-scale fine-grained building facade point clouds dataset, which can be used to facilitate research on visual tasks related to building facade. In order to fully explore the local neighborhood features in the 3D point cloud on semantic segmentation task, we present a novel attention-based network DLA-net in which the self-attention block and attentive pooling block based on the powerful attention mechanism are used for learning important local features. DLA module could be easily embedded into various network architectures for point cloud segmentation and we embed it into an encoder-decoder architecture, resulting in the DLA-Net in this work. Extensive experiments on the building facade dataset benchmarks demonstrate the state-of-the-art performance of our proposed DLA-Net. In the future, the dataset can be used not only for the task of semantic segmentation, but also for many computer vision tasks such as building facade point clouds completion and 3D reconstruction.

## References

[1] Segmentation and reconstruction of polyhedral building roofs from aerial lidar point clouds, IEEE Transactions on Geoscience and Remote Sensing 48 (3 PART2) (2010) 1554–1567.

[2] Layered analysis of irregular facades via symmetry maximization, ACM Transactions on Graphics 32 (4) (2013).

[3] Tensor discriminative locality alignment for hyperspectral image spectral-spatial feature extraction, IEEE Transactions on Geoscience and Remote Sensing 51 (1) (2013) 242–256.

[4] M. Zhang, L. Zhang, P. Takis Mathiopoulos, Y. Ding, H. Wang, Perception-based shape retrieval for 3D building models, ISPRS Journal of Photogrammetry and Remote Sensing 75 (2013) 76–91.

[5] K. R. Nia, G. Mori, Building damage assessment using deep learning and ground-level image data, in: 14th Conference on Computer and Robot Vision (CRV), 2017, pp. 95–102.

[6] C. Yang, T. Han, L. Quan, C.-L. Tai, Parsing façade with rank-one approximation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2012, pp. 1720–1727.

[7] D. Dai, M. Prasad, G. Schmitt, L. Van Gool, Learning domain knowledge for facade labelling, in: Proceedings of the European Conference on Computer Vision (ECCV), 2012, pp. 710–723.

[8] A. Martinović, M. Mathias, J. Weissenberg, L. Van Gool, A three-layered approach to facade parsing, in: Proceedings of the European Conference on Computer Vision (ECCV), 2012, pp. 416–429.

[9] R. Tyleček, R. Šára, Spatial pattern templates for recognition of objects with regular structure, in: Proceedings of German Conference on Pattern Recognition (GCPR), 2013, pp. 364–374.

[10] H. Riemenschneider, A. Bódis-Szomorú, J. Weissenberg, L. Van Gool, Learning where to classify in multi-view semantic segmentation, in: Proceedings of the European Conference on Computer Vision (ECCV), 2014, pp. 516–532.

[11] A. Martinovic, J. Knopp, H. Riemenschneider, L. Van Gool, 3D all the way: Semantic segmentation of urban scenes from start to end in 3D, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015, pp. 4456–4465.

[12] A. Serna, B. Marcotegui, F. Goulette, J.-E. Deschaud, Paris-rue-Madame database: a 3D mobile laser scanner dataset for benchmarking urban detection, segmentation and classification methods, in: 4th International Conference on Pattern Recognition, Applications and Methods (ICPRAM), 2014.

[13] X. Roynard, J.-E. Deschaud, F. Goulette, Paris-Lille-3D: A large and high-quality ground-truth urban point cloud dataset for automatic segmentation and classification, The International Journal of Robotics Research 37 (6) (2018) 545–557.

[14] Y. Guo, H. Wang, Q. Hu, H. Liu, L. Liu, M. Bennamoun, Deep Learning for 3D Point Clouds: A Survey, IEEE Transactions on Pattern Analysis and Machine Intelligence (2020) 1–1.

[15] N. Audebert, B. Le Saux, S. Lefèvre, Semantic segmentation of earth observation data using multimodal and multi-scale deep networks, in: Proceedings of Asian Conference on Computer Vision (ACCV), 2016, pp. 180–196.

[16] F. J. Lawin, M. Danelljan, P. Tosteberg, G. Bhat, F. S. Khan, M. Felsberg, Deep projective 3D semantic segmentation, in: International Conference on Computer Analysis of Images and Patterns (CAIP), 2017, pp. 95–107.

[17] M. Tatarchenko, J. Park, V. Koltun, Q.-Y. Zhou, Tangent convolutions for dense prediction in 3D, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 3887–3896.

[18] Y. Lin, Z. Yan, H. Huang, D. Du, L. Liu, S. Cui, X. Han, FPconv: Learning local flattening for point convolution, in: Proceedings of the IEEE

Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 4293–4302.

[19] J. Long, E. Shelhamer, T. Darrell, Fully convolutional networks for semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2015, pp. 3431–3440.

[20] J. Huang, S. You, Point cloud labeling using 3D convolutional neural network, in: 23rd International Conference on Pattern Recognition (ICPR), 2016, pp. 2670–2675.

[21] L. Tchapmi, C. Choy, I. Armeni, J. Gwak, S. Savarese, Segcloud: Semantic segmentation of 3D point clouds, in: International conference on 3D Vision (3DV), 2017, pp. 537–547.

[22] A. Dai, D. Ritchie, M. Bokeloh, S. Reed, J. Sturm, M. Nießner, ScanComplete: Large-scale scene completion and semantic segmentation for 3D scans, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 4578–4587.

[23] C. R. Qi, H. Su, K. Mo, 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 (CVPR), 2017, pp. 652–660.

[24] C. R. Qi, L. Yi, H. Su, L. J. Guibas, PointNet++: Deep hierarchical feature learning on point sets in a metric space, in: Advances in Neural Information Processing Systems (NeurIPS), 2017, pp. 5099–5108.

[25] H. Zhao, L. Jiang, C.-W. Fu, J. Jia, PointWeb: Enhancing local neighborhood features for point cloud processing, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 5565–5573.

[26] Q. Hu, B. Yang, L. Xie, S. Rosa, Y. Guo, Z. Wang, N. Trigoni, A. Markham, RandLA-Net: Efficient Semantic Segmentation of Large-

Scale Point Clouds, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 11108–11117.

[27] Y. Li, R. Bu, M. Sun, W. Wu, X. Di, B. Chen, PointCNN: Convolution on x-transformed points, in: Advances in Neural Information Processing Systems (NeurIPS), 2018, pp. 820–830.

[28] M. Jiang, Y. Wu, T. Zhao, Z. Zhao, C. Lu, PointSIFT: A SIFT-like Network Module for 3D Point Cloud Semantic Segmentation, arXiv preprint arXiv:1807.00652 (2018).

[29] H. Thomas, C. R. Qi, J.-E. Deschaud, B. Marcotegui, F. Goulette, L. J. Guibas, KPConv: Flexible and Deformable Convolution for Point Clouds, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 6411–6420.

[30] X. Wang, J. He, L. Ma, Exploiting local and global structure for point cloud semantic segmentation with contextual point representations, in: Advances in Neural Information Processing Systems (NeurIPS), 2019.

[31] C. Zhao, W. Zhou, L. Lu, Q. Zhao, Pooling scores of neighboring points for improved 3D point cloud segmentation, in: Proceedings of IEEE International Conference on Image Processing (ICIP), 2019, pp. 1475–1479.

[32] L. Wang, Y. Huang, Y. Hou, S. Zhang, J. Shan, Graph attention convolution for point cloud semantic segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 10288–10297.

[33] M. Feng, L. Zhang, X. Lin, S. Z. Gilani, A. Mian, Point attention network for semantic segmentation of 3D point clouds, Pattern Recognition 107 (2020) 107446.

[34] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin, Attention is all you need, in: Advances in Neural Information Processing Systems (NeurIPS), 2017.

[35] F. Wu, A. Fan, A. Baevski, Y. N. Dauphin, M. Auli, Pay less attention with lightweight and dynamic convolutions, in: International Conference on Learning Representations (ICLR), 2019, pp. 1–14.

[36] H. Hu, Z. Zhang, Z. Xie, S. Lin, Local relation networks for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 3464–3473.

[37] P. Ramachandran, N. Parmar, A. Vaswani, I. Bello, A. Levskaya, J. Shlens, Stand-alone self-attention in vision models, in: Advances in Neural Information Processing Systems (NeurIPS), 2019, pp. 68–80.

[38] H. Zhao, J. Jia, V. Koltun, Exploring self-attention for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 10076–10085.

[39] A. Boulch, B. Le Saux, N. Audebert, Unstructured point cloud semantic labeling using deep segmentation networks., Eurographics Workshop on 3D Object Retrieval (2017) 17–24.

[40] T. Le, Y. Duan, Pointgrid: A deep network for 3D shape understanding, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 9204–9214.

[41] H.-Y. Meng, L. Gao, Y.-K. Lai, D. Manocha, Vv-net: Voxel vae net with group convolutions for point cloud segmentation, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 8500–8508.

[42] D. Rethage, J. Wald, J. Sturm, N. Navab, F. Tombari, Fully-convolutional point networks for large-scale point clouds, in: Proceedings of the European Conference on Computer Vision (ECCV), 2018, pp. 596–611.

[43] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, J. M. Solomon, Dynamic graph cnn for learning on point clouds, ACM Transactions on Graphics 38 (5) (2018).

[44] Q. Huang, W. Wang, U. Neumann, Recurrent Slice Networks for 3D Segmentation of Point Clouds, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 2626–2635.

[45] B. Yang, S. Wang, A. Markham, N. Trigoni, Robust attentional aggregation of deep feature sets for multi-view 3D reconstruction, International Journal of Computer Vision 128 (1) (2020) 53–73.

[46] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2016, pp. 770–778.

[47] RIEGL VMX-450 Datasheet. Accessed 2-May-2015, `http://www.riegl.com/nc/products/mobile-scanning/produktdetail/product/scannersystem/10/`.