# PointeNet: A Lightweight Framework for Effective and Efficient Point Cloud Analysis

Lipeng Gu[a], Xuefeng Yan[a], Liangliang Nan[b], Dingkun Zhu[c], Honghua Chen[a], Weiming Wang[d], Mingqiang Wei[a]

[a]*Nanjing University of Aeronautics and Astronautics*
[b]*Delft University of Technology*
[c]*Jiangsu University of Technology*
[d]*Hong Kong Metropolitan University*

## Abstract

Current methodologies in point cloud analysis predominantly explore 3D geometries, often achieved through the introduction of intricate learnable geometric extractors in the encoder or by deepening networks with repeated blocks. However, these approaches inevitably lead to a significant number of learnable parameters, resulting in substantial computational costs and imposing memory burdens on CPU/GPU. Additionally, the existing strategies are primarily tailored for object-level point cloud classification and segmentation tasks, with limited extensions to crucial scene-level applications, such as autonomous driving. In response to these limitations, we introduce PointeNet, an **e**fficient **net**work designed specifically for **point** cloud analysis. PointeNet distinguishes itself with its lightweight architecture, low training cost, and plug-and-play capability, effectively capturing representative features. The network consists of a Multivariate Geometric Encoding (MGE) module and an *optional* Distance-aware Semantic Enhancement (DSE) module. The MGE module employs operations of sampling, grouping, and multivariate geometric aggregation to lightweightly capture and adaptively aggregate multivariate geometric features, providing a comprehensive depiction of 3D geometries. The DSE module, designed for real-world autonomous driving scenarios, enhances the semantic perception of point clouds, particularly for distant points. Our method demonstrates flexibility by seamlessly integrating with a classification/segmentation head or embedding into off-the-shelf 3D object detection networks, achieving notable performance improvements at a minimal cost. Extensive experiments on object-level datasets, including ModelNet40, ScanObjectNN, ShapeNetPart, and the scene-level dataset KITTI, demonstrate the superior performance of PointeNet over state-of-the-art methods in point cloud analysis. Notably, PointeNet outperforms PointMLP with significantly fewer parameters on ModelNet40, ScanObjectNN, and ShapeNetPart, and achieves a substantial improvement of nearly 3% in 3D $AP_{R40}$ for PointRCNN on KITTI with a minimal parameter cost of 1.4 million.

*Keywords:* PointeNet, efficient point cloud analysis, lightweight framework

## 1. Introduction

With the popularity of 3D sensors such as LiDARs and depth cameras Geiger et al. (2012), point cloud analysis has gained significant prominence in both academic research and industrial development Qi et al. (2017a,b, 2019); Liang et al. (2022). Unlike grid-based RGB images, point clouds comprise unordered and irregular points that outline the surfaces of objects in 3D space. This characteristic poses significant challenges in designing effective point cloud analysis models.

PointNet/PointNet++ Qi et al. (2017a,b) are pioneering works capable of directly analyzing unordered point clouds without the need for preprocessing. Just as ResNet He et al. (2016) serves as a prominent backbone in image processing, PointNet/PointNet++ has emerged as a widely adopted backbone network in subsequent point cloud analysis models Qi et al. (2019); Yang et al. (2020); Shi et al. (2019); Ma et al. (2022). The iterative evolution in these methods adheres to the "*continual increment*" principle, primarily focusing on modifying the backbone network, PointNet++, to capture more representative local geometric features. Specific improvements fall mainly into two aspects: introducing well-designed local geometric extractors (such as Graph Convolution Wang et al. (2019); Lin et al. (2022), Adaptive Point Convolution Xu et al. (2021), and Residual Point Block Ma et al. (2022)) into the encoder or simply stacking repeated blocks to deepen the network Ma et al. (2022); Liu et al. (2019). Such strategies, while capable of boosting performance, inevitably lead to a notable surge in network parameters. Fortunately, Point-NN recognizes this longstanding challenge and for the first time attempts to employ a counterintuitive "*subtraction*" strategy by trimming PointNet++, retaining only its non-learnable components—sampling, grouping, and pooling. This allows tasks like point cloud classification, segmentation, and even 3D detection to be performed without the need for training. While this unique strategy significantly reduces training costs and achieves point cloud analysis, it falls short of exploiting the local geometry properties of point clouds, relying solely on spatial neighboring information. Furthermore, existing methods, including Point-NN, primarily focus on object-level point cloud analysis and see a limited extension to more valuable scene-level applications, such as autonomous driving. These problems hinder Point-NN from fully unleashing its potential.

In this paper, we propose an e**ffi**cient **Net**work for ***point*** cloud analysis, dubbed ***PointeNet***. Our model is highly streamlined, comprising only non-parametric computational components such as sampling, grouping, pooling, and a minimal number of learnable parameters, and can be flexibly combined with a point cloud classification/segmentation head, or be embedded into cutting-edge 3D detection networks tailored for real-world autonomous driving scenarios to enhance performance. Specifically, PointeNet comprises two crucial modules: the multivariate geometric encoding (MGE) module, and the distance-aware semantic enhancement (DSE) module. MGE comprises non-learnable farthest point sampling (FPS) and k-nearest neighbor (k-NN) components, along with a multivariate local geometric aggregation (MLGA) module incorporating a minimal number of learnable parameters. This module captures multivariate 3D geometric features within local regions of point clouds, encompassing curvature, normal, and spatial neighboring information. DSE is an optional module tailored for autonomous driving scenarios. It dynamically adjusts the segmentation difficulty on a point-by-point basis, particularly allocating more "attention" to challenging distant points, and outputs distance-aware semantic features to rich point clouds, further enhancing the performance of arbitrary 3D detection networks. We evaluate the performance of PointeNet through a comparative analysis with nineteen competitors across the ScanObjectNN Uy et al. (2019), ModelNet40 Wu et al. (2015), ShapeNetPart Yi et al. (2016), and KITTI Geiger et al. (2012) datasets. Remarkably, PointeNet surpasses all competitors, showcasing its superior performance. The contributions of this work are threefold:

- We analyze the longstanding challenge in recent point cloud analysis methods and propose an efficient network, dubbed PointeNet. It excels in lightweight capturing and adaptive aggregation of multivariate geometric and semantic features, making it suitable for point cloud segmentation/classification and enhancing arbitrary 3D object detection networks tailored for autonomous driving scenarios.

- We propose a multivariate geometric encoding (MGE), featuring non-learnable FPS and k-NN modules, complemented by a minimal number of learnable parameters, to capture diverse 3D geometric features.

- We propose an optional distance-aware semantic enhancement (DSE) module tailored for autonomous driving scenarios. It dynamically adjusts the segmentation difficulty on a point-by-point basis, giving more attention to challenging distant points, and outputs distance-aware semantic features.

## 2. Related Work

### 2.1. Point Cloud Analysis

There are two main paradigms for processing irregular and unordered point clouds: grid-based and point-based methods. Grid-based methods Liang et al. (2022); Qi et al. (2019); Shi et al. (2020); Lang et al. (2019) involve projecting irregular point clouds onto regular grids, such as pillars or voxels, and then processing them using 2D/3D convolutional neural networks (CNNs). This paradigm significantly enhances the processing speed by transforming the point cloud into structured grids. However, projecting onto grids may lead to information loss, degrading the quality of the point cloud representation Yang et al. (2019).

In contrast, point-based methods have emerged to directly process irregular point clouds without additional regularization, preserving the point cloud details more accurately. PointNet Qi et al. (2017a) is a pioneering method that utilizes shared MLPs to handle unordered point clouds as input directly. Point-Net++ Qi et al. (2017b) builds upon PointNet by introducing a hierarchical feature learning paradigm, allowing for the recursive capture of local geometric structures. Due to the promising performance exhibited by PointNet++, particularly in leveraging local geometric representations, including multi-scale geometric information, it has become a foundational element in modern point cloud analysis methods. Subsequent works, such as KPConv Thomas et al. (2019), RSCNN Liu et al. (2019), 3D-GCN Lin et al. (2022), PAConv Xu et al. (2021), PointConv Wu et al. (2019) and PointMLP Ma et al. (2022), either introduce intricate local geometric extractors or simply stack repeated network blocks to achieve further performance improvements. However, these strategies make the networks more complex and less efficient. *In contrast, our PointeNet draws from the successful experience of Point-NN and proposes an efficient network primarily based on non-learnable components, capturing representative multivariate geometric and semantic features at a minimal cost in terms of parameters.*

### 2.2. Local Geometry Exploration

Since the powerful ability to capture local geometric features of point clouds demonstrated by Point-Net++ Qi et al. (2017b), subsequent research has mainly focused on exploring local geometric representations, which can be divided into three categories: convolution-based, graph-based, and attention-based methods. Classic convolution-based methods, with PointConv Wu et al. (2019) as a representative example, employ MLPs to approximate continuous weights and density functions within convolution filters. It extends dynamic filtering to a novel convolution operation. Unlike convolution-based methods, graph-based methods explore relationships between points. For instance, DGCNN Wang et al. (2019) introduces EdgeConv, a novel method that generates edge features describing the relationships between points and their neighbors, capturing local features of the point cloud. 3D-GCN Lin et al. (2022) employs a 3D graph convolutional network to form deformable 3D kernels, directly performing convolutional computations on point clouds. Related to graph-based methods, attention-based methods similarly focus on exploring direct relationships between points, as seen in works like PCT Guo et al. (2021) and Point Transformer Zhao et al. (2021). Although these methods have achieved success in point cloud analysis, they have consistently employed an "*increment*" principle, undoubtedly introducing more learnable parameters and thereby reducing computational efficiency and exacerbating the burden on CPU/GPU. Compared with these methods, our PointeNet achieves better performance with fewer parameters.

### 2.3. Deep Network Architecture for Point Clouds

The advancement of methods in point cloud analysis often coincides with breakthroughs in image processing networks. For instance, following the success of ResNet He et al. (2016) based on simple convolutions in the field of image processing, subsequent methods often use it as a backbone and continuously improve upon it. Similarly, after the breakthrough achieved by PointNet/PointNet++ based on simple MLPs in point cloud processing, they have become the mainstream backbones for point cloud analysis. After the success of graph- Felzenszwalb and Huttenlocher (2004), attention- Wang et al. (2018), and transformer-based Dosovitskiy et al. (2021) methods in the field of image processing, they have inspired a series of works in point cloud analysis, contributing to the further development of point cloud analysis methods Wu et al.
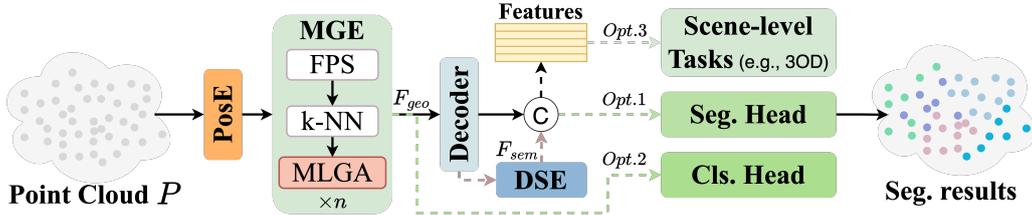
Figure 1: **Overview of PointeNet.** It is an efficient network comprising two key modules: the multivariate geometric encoding (MGE) module and the distance-aware semantic enhancement (DSE) module. It can be flexibly combined with segmentation/classification heads to achieve excellent point cloud segmentation/classification tasks. Furthermore, it can serve as a plug-and-play module embedded in arbitrary scene-level tasks, such as 3D Object Detection (3OD), contributing to further performance enhancement.

(2019); Wang et al. (2019); Ran et al. (2021); Guo et al. (2021); Zhao et al. (2021). Our PointeNet incorporates the strengths of previous methods, enabling the lightweight capturing and adaptive aggregation of diverse geometric and semantic representations of point clouds at a lower cost, thereby achieving powerful point cloud analysis capabilities.

## 3. Methodology

### 3.1. Overview

We propose an efficient network for point cloud analysis, named PointeNet. It can lightweightly capture and adaptively aggregate multivariate geometric and semantic features of point clouds using a minimal number of learnable parameters. The detailed framework of PointeNet is illustrated in Fig. 1.

### 3.2. Revisiting Point-based Methods

Unlike grid-based methods, point-based methods have gained popularity due to their ability to directly learn the underlying point cloud representation without preprocessing. PointNet/PointNet++ Qi et al. (2017a,b) are pioneers in this paradigm, employing a hierarchical feature learning approach through the stacking of multiple learning stages. We first delve into the core idea of PointNet++ and then explore the strengths and weaknesses of subsequent methods.

Given a set of points $\mathcal{P} = \{p_i \mid i = 1, \ldots, \mathcal{N}\} \in \mathbb{R}^{\mathcal{N} \times 3}$, where $\mathcal{N}$ is the number of points. The core idea of PointNet++ is to utilize FPS for down-sampling point clouds, use k-NN for grouping points, and then employ local geometric extractors to extract features, followed by max pooling to aggregate these features. The key steps for geometric feature extraction and aggregation can be formulated as:

$$f_c = \mathcal{A}\left(\Phi\left(f_{c,j}\right) \mid j \in \mathcal{N}_c\right) \tag{1}$$

where $\mathcal{A}(\cdot)$ means aggregation function (i.e., max-pooling in PointNet++), $\Phi(\cdot)$ denotes the local feature extraction extractor (i.e., MLPs in PointNet++), and $f_{c,j}$ is the $j-$th neighbor point feature of center point $p_c$ within the local region.

Regarding network architecture design, PointNet++ establishes a versatile baseline for point cloud analysis. Later, based on this baseline, several methods are introduced, either emphasizing local geometric extractors Wang et al. (2019); Lin et al. (2022); Thomas et al. (2019); Xu et al. (2021) or stacking repeated network blocks to deepen the network Ma et al. (2022); Liu et al. (2019). While these methods can capture detailed local geometric information and often yield promising performance, their development encounters three challenges.

Firstly, whether by introducing intricate geometric extractors or deepening the network to enhance the representativeness of local geometric features, these methods inevitably introduce a large number of learnable parameters, leading to increased computational complexity and exacerbating the burden on CPU/GPU,

4

resulting in inference latency. For instance, PointMLP Ma et al. (2022) achieves 16.7 M parameters by stacking residual network blocks, and the training time is as high as 23 hours when performing segmentation tasks on ShapeNetPart. This represents a considerable computational overhead. Secondly, intricate local geometric extractors are maturing and even coming to saturation. Performance improvements on popular benchmarks like ModelNet40, ScanObjectNN, and ShapeNetPart are encountering bottlenecks. Third, these methods are primarily designed for object-level point cloud analysis scenarios and rarely consider more realistic and valuable scene-level applications, such as autonomous driving. These limitations prompt us to explore a new paradigm that "lightens the load" on local feature extractors. In other words, *can we design an efficient lightweight point cloud analysis framework with fewer parameters that achieves satisfactory performance not only at the object level but also at the scene level?*

Excitingly, Point-NN makes the first successful attempt to perform point cloud analysis by constructing a PointNet++-style hierarchical framework using only non-learnable components such as FPS, k-NN, and pooling. However, it relies solely on a single geometric feature (i.e., spatial neighborhood information), overlooking the more potential point cloud features, such as normals and curvatures of the local surfaces, and even semantic features. And, it also lacks customized optimization for real-world scene-level applications.

### 3.3. Framework of PointeNet

To overcome the aforementioned limitations, we propose a simple yet efficient network, named PointeNet. It can lightweightly capture and adaptively aggregate multivariate geometric and semantic features without the need for intricate or heavy operations. Moreover, it is specifically optimized for real-world autonomous driving scenarios.

As shown in Figure 1, the key operations of PointeNet can be formulated as:

$$f_c = \mathcal{C}\left(f_c^g, \Phi_{sem}\left(f_c^g\right)\right) \tag{2}$$

$$f_c^g = \mathcal{A}\left(\left(\Phi_{pos}\left(f_{c,j}\right), \Phi_{sur}\left(f_{c,j}\right)\right) \mid j \in \mathcal{N}_c\right) \tag{3}$$

where $f_{c,j}$ is the $j-$th neighbor point feature of the center point $p_c$ within the local region. $\Phi_{pos}\left(\cdot\right)$ and $\Phi_{sur}\left(\cdot\right)$ are functions in the multivariate geometric encoding module responsible for extracting spatial neighboring features and curvature-normal features, respectively. And, $\Phi_{sem}\left(\cdot\right)$ is a function for extracting distance-aware semantic features tailored for real-world autonomous driving scenarios. $\mathcal{A}$ is an aggregation function, and in PointeNet, it is our proposed multivariate adaptive aggregation (MAA) module. $\mathcal{C}$ is the concatenation operation, which concatenates the multivariate geometric features and distance-aware semantic features.

Our efficient PointeNet exhibits some prominent advantages: i) PointeNet comprises only a minimal number of parameters for positional encoding, feature aggregation, and semantic segmentation. It is naturally unaffected by permutations, aligning perfectly with the characteristics of point clouds. ii) Without the need for designing and stacking intricate feature extractors, PointeNet efficiently achieves state-of-the-art performance with only a minimal number of parameters. iii) Additionally, due to the absence of intricate feature extractors, PointeNet is easy to train and has low hardware requirements.

### 3.4. Multivariate Geometric Encoding

To clearly describe a given point cloud, our natural intuition leads us to seek mastery over its spatial neighboring, surface detail, and even semantic information. However, existing methods,



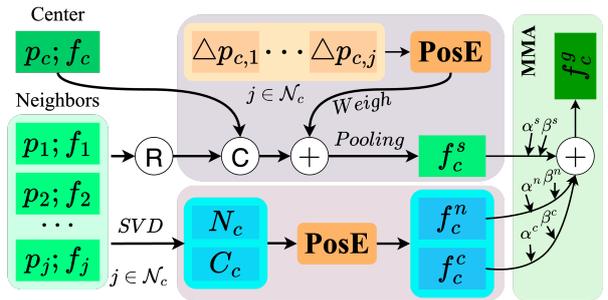Figure 2: **Overview of multivariate geometric aggregation (MLGA) module.** It is capable of lightweight capturing of 3D geometries, including the spatial neighboring $f_c^s$, curvature $f_c^c$, and normal $f_c^n$ features within a local region, and then undergoes a multivariate adaptive aggregation (MAA) module to adaptively Aggregate various features, finally outputting multivariate local geometric features $f_c^g$

5

such as PointMLP Ma et al. (2022) and Point-NN Zhang et al. (2023), almost entirely overlook the potential surface details or employ elaborate fusion strategies to aggregate various features. In this regard, we introduce a Multivariate Geometric Encoding (MGE) module, consisting of non-parametric FPS, k-NN, and a module for multivariate local geometric aggregation (MLGA). This module is designed to efficiently capture and adaptively aggregate multivariate geometric features.

MLGA is a crucial component of MGE. As illustrated in Figure 2, it not only captures spatial neighboring information within the local region but also comprehensively depicts local surface details, including curvature and normal information. Moreover, it achieves adaptive aggregation of multivariate geometric features. Note that a limited number of parameters in MLGA exist in the position encoding and the adaptive feature aggregation process. In the following, we elaborate on the process of capturing and aggregating multivariate geometric features.

The encoding of spatial neighboring information in local neighborhoods, similar to Point-NN, is represented by the following formula:

$$f_c^s = \left( \mathcal{C} \left( f_c, f_j \right) + PosE \left( \triangle p_{c,j} \right) \right) \odot PosE \left( \triangle p_{c,j} \right), j \in \mathcal{N}_c \tag{4}$$

where $f_c$ and $f_j$ represent the center and neighbor features within the local region, respectively. $\triangle p_{c,j}$ represents the normalized coordinates of neighboring points by the mean and standard deviation. $PosE\left( \cdot \right)$ refers to the position encoding, and here, it is implemented using an MLPs.

Surface feature information in point clouds refers to a set of features that represent the shape of surfaces in 3D space. These features are crucial for point cloud analysis, and common surface feature information includes normals, curvature, smoothness, boundaries, and even texture. Among these, normals and curvature are the most typical surface features. GeoMAE has successfully demonstrated this by using normals and curvature of local surfaces in point clouds as self-supervised signals for pre-training, leading to improved performance in downstream tasks such as 3D detection. Therefore, to enhance the representation of the single spatial neighboring features in the encoder, we introduce surface curvature and normal as two additional geometries. For a set of neighboring points $p_j$, we begin by computing the covariance matrix:

$$M_c = \frac{1}{\mathcal{N}_c} \sum_{j=1}^{\mathcal{N}_c} p_j p_j^T - \bar{p}\bar{p}^T \tag{5}$$

where $\bar{p} = \frac{1}{\mathcal{N}_c} \sum_{j=1}^{\mathcal{N}_c} p_j$ is the centroid of the local region. Next, we perform Singular Value Decomposition (SVD) on the covariance matrix $M_c$ to obtain the eigenvalues $c_1, c_2, c_3$ and corresponding eigenvectors $n_1, n_2, n_3$. Here, we use the eigenvector $n_3$ corresponding to the smallest eigenvalue $\lambda_3$ as the pseudo-normal vector $N_c = n_3$. Following that, we normalize the three eigenvalues to obtain the pseudo-curvature vector $C_c = \{c_0, c_1, c_2\}$:

$$c_m = \frac{c_m}{\sum_{i=1}^3 c_i}, m \in \{1, 2, 3\} \tag{6}$$

As both pseudo-curvature and pseudo-normal vectors are low-dimensional, to embed them into high-dimensional geometric features, we naturally use positional encoding to map them into a higher-dimensional space, i.e., $f_c^n = PosE\left( N_c \right)$ and $f_c^c = PosE\left( C_c \right)$.

Unlike previous methods Zhang et al. (2019) that use MLPs with relatively more parameters to aggregate multiple features, we ingeniously introduce learnable parameters $\alpha, \beta$ to learn channel-wise weights and biases for each feature, i.e., $f_c^s$, $f_c^n$ and $f_c^c$. We then perform adaptive aggregation using the following formula:

$$f_c^g = \sum_{i=\{s,c,n\}} \left( \alpha^i \odot f_c^i + \beta^i \right) \tag{7}$$

where $\odot$ denotes Hadamard product and $f_c^g$ is the output of aggregated multivariate geometric features.

### 3.5. Distance-ware Semantic Enhancement

Prior point cloud analysis methods, such as PointMLP Ma et al. (2022) and Point-NN Zhang et al. (2023), are previously limited to object-level classification/segmentation tasks and do not have custom optimizations for scene-level tasks, such as applications in autonomous driving scenarios. In real-world autonomous driving applications, the field of view is more expansive, and the point cloud becomes sparser, particularly as the distance increases. This poses a challenge for effective perception, especially for distant points.

The intuitive solution is to allocate more "attention" to the sparser point clouds at a distance, while the relatively denser point clouds nearby require only a small amount of "attention" for effective perception. PV-RCNN Shi et al. (2020) attempts to give more attention to foreground points through semantic segmentation. However, it neglects the varying segmentation difficulties in different regions of the entire point cloud scene based on their distances. Therefore, the key is to distribute attention from the denser point clouds nearby to the sparser regions at a distance, ensuring a more equitable and easily trainable perception network.

To address this, we introduce the Distance-aware Segmentation Enhancement (DSE) module (see Figure 3), which incorporates the distance factor into the semantic segmentation loss, specifically using Focal Loss Lin et al. (2020). DSE takes the multivariate geometric features $f^g$ from MGE and corresponding point coordinates $p$ as input and outputs the features of the intermediate segmentation process as the semantic features $f^s$. Specifically, DSE employs four fully connected (FC) layers and the *Sigmoid* function to perform binary classification for each point (foreground or background). Additionally, it utilizes two FC layers and the *Softmax* function to extract the distance factor $d \in [0,1]$ based on the absolute values of the $x, y$ coordinates for each point $p$. Subsequently, we incorporate $d$ into Focal Loss, replacing the original adjustment factors $\alpha$ and $\gamma$ to adaptively adjust the segmentation penalty based on the distance, as expressed by the formula:

$$L_{seg} = -d\,(1 - p_t)^{\frac{1}{d}} \log p_t \tag{8}$$

where $p_t$ denotes the probability of the predicted sample belonging to the positive class, the term $-d\,(1 - p_t)^{\frac{1}{d}}$ serves as the weight for the $\log p_t$ term, with $d$ acting as an adaptive adjustment factor for this weight. For distant points, their $x, y$ coordinates are relatively large, leading to larger distance factors $d$. Conversely, for nearby points, their distance factors $d$ are smaller. Consequently, $\frac{1}{d}$ for distant points is smaller than for nearby points, resulting in a greater penalty for distant points represented by $(1 - p)^{\frac{1}{d}} \log p_t$ compared to nearby points. Furthermore, as the distance factor $d$ increases, the penalty $(1 - p)^{\frac{1}{d}} \log p_t$ also increases. To regulate the overall loss value, we multiply it by the distance factor $d$.

By directly concatenating multivariate geometric features and distance-aware semantic features, we utilize the combined features as auxiliary features for point clouds. This allows for a cost-effective enhancement of the performance of any 3D object detection network. This improvement not only requires minimal code modifications but also incurs low computational costs. *Importantly, it enables the extension of point cloud analysis methods to more practically valuable scene-level tasks.*



Figure 3: **Overview of distance-ware semantic enhancement (DSE) module.** It introduces the distance factor $d$ into the semantic segmentation loss function, adaptively distinguishing the segmentation difficulty in different regions to obtain distance-aware semantic features $f_c^s$.

## 4. Experiments

In this section, we conduct a comprehensive evaluation of PointeNet not only on object-level datasets such as ModelNet40 Wu et al. (2015), ScanObjectNN Uy et al. (2019), and ShapeNetPart Yi et al. (2016), but also validate its performance on a scene-level dataset KITTI Geiger et al. (2012). Detailed ablation studies demonstrate the effectiveness and efficiency of PointeNet with both quantitative and qualitative analysis.
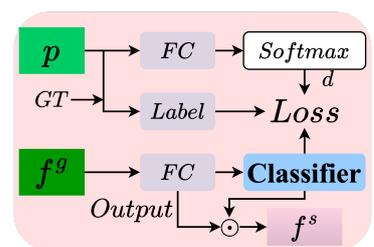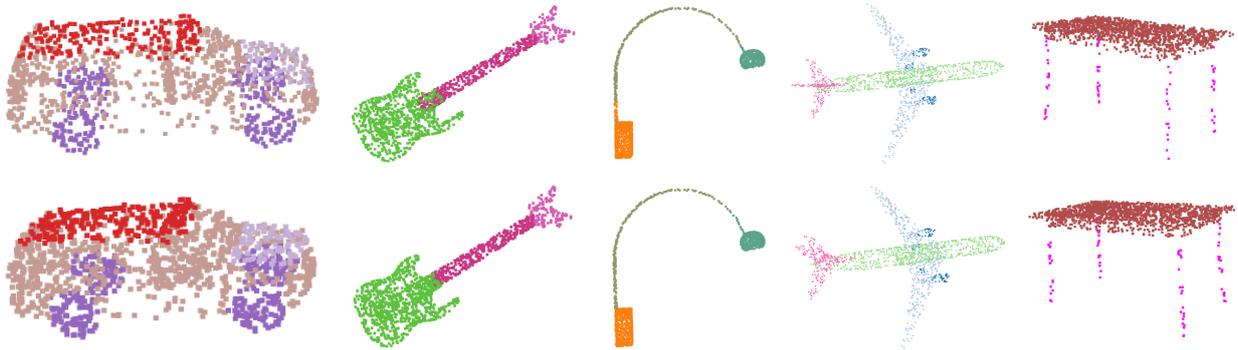
Figure 4: **Visualization of part segmentation results on ShapeNetPart**. The top line represents the ground truth, while the bottom line corresponds to our predictions. Upon visualization, our predictions closely align with the ground truth.

### 4.1. Datasets

For *Part Segmentation*, we employ the ShapeNetPart Yi et al. (2016) dataset, which consists of 16881 shapes with 16 classes, totaling 50 part labels. In each class, the number of parts varies from 2 to 6.

For *Shape Classification*, we evaluate the performance on the synthetic ModelNet40 Wu et al. (2015) and the real-world ScanObjectNN Uy et al. (2019) datasets. ModelNet40 comprises 9843 training and 2468 testing meshed CAD models distributed across 40 categories. ScanObjectNN includes 15000 objects categorized into 15 classes, featuring 2902 unique object instances in real-world scenarios. Due to the presence of background, noise, and occlusions, ScanObjectNN poses significant challenges for existing methods. In our experiments, we specifically consider the most challenging perturbed variant.

For *3D Object Detection*, our experiments are conducted on the KITTI dataset Geiger et al. (2012), which is a popular benchmark for 3OD in autonomous driving scenes. It contains 7481 samples for training and 7518 samples for testing. Each sample consists of a point cloud and an RGB image with nine categories.

Table 1: **Part segmentation on ShapeNetPart.** We report mean IoU scores across classes (cls. mIoU) and instances (Inst. mIoU) on the testing set. Note that the embedding dimension of PointeNet is set to 90 and ∗ means the reproduced results.

| Method | Cls. mIoU (%) | Inst. mIoU (%) | Param. | Train Time | Test Time |
|---|---|---|---|---|---|
| PointNet Qi et al. (2017a) | 80.4 | 83.7 | 8.3 M | - | - |
| PointNet++ Qi et al. (2017b) | 81.9 | 85.1 | 1.8 M | - | - |
| Kd-Net Klokov and Lempitsky (2017) | - | 82.3 | - | - | - |
| SpiderCNN Xu et al. (2018) | 82.4 | 85.3 | - | - | - |
| SPLATNet Su et al. (2018) | 83.7 | 85.4 | - | - | - |
| PointMLP∗ Ma et al. (2022) | 84.0 | **85.7** | 16.7 M | 23 h | 7 ms |
| Point-NN∗ Zhang et al. (2023) | - | 70.7 | 0 M | 0 h | 19 ms |
| PointeNet (Ours) | **84.2** | 85.6 | 8.6 M | 8 h | 4 ms |

### 4.2. Part Segmentation on ShapeNetPart

We present the results of PointeNet and compare our method with several recent works, including PointNet Qi et al. (2017a), PointNet++ Qi et al. (2017b), Kd-Net Klokov and Lempitsky (2017), SpiderCNN Xu et al. (2018), SPLATNet Su et al. (2018), PointMLP Ma et al. (2022) and Point-NN Zhang et al. (2023), for the 3D shape part segmentation task on the ShapeNetPart benchmark, as summarized in Table 1.

Analyzing the results, our PointeNet achieves segmentation performance comparable to PointMLP while utilizing only half of its parameters (8.6 M vs. 16.7 M), especially surpassing PointMLP by 0.2% in cls. mIoU (84.2 % vs. 84.0 %). Additionally, our PointeNet exhibits a clear advantage over PointMLP in terms

Table 2: **Shape classification on synthetic ModelNet40.** We report the class-average accuracy (mAcc) and overall accuracy (OA) on the testing set. ∗ means the reproduced results.

| Method | mAcc (%) | OA (%) | Param. | Train Time | Test Time |
|---|---|---|---|---|---|
| PointNet Qi et al. (2017a) | 86.0 | 89.2 | - | - | - |
| PointNet++ Qi et al. (2017b) | - | 91.9 | 1.4 M | - | - |
| PointCNN Li et al. (2018) | 88.1 | 92.5 | - | - | - |
| PointConv Wu et al. (2019) | - | 92.5 | 18.6 M | - | - |
| KPConv Thomas et al. (2019) | - | 92.9 | 15.2 M | - | - |
| Point Trans. Zhao et al. (2021) | 90.6 | 93.7 | - | - | - |
| GBNet Qiu et al. (2022) | 91.0 | 93.8 | 8.4 M | - | - |
| PointMLP∗ Ma et al. (2022) | 90.8 | 93.3 | 13.2 M | 11 h | 5 ms |
| Point-NN∗ Zhang et al. (2023) | - | 81.3 | 0 M | 0 h | 4 ms |
| PointeNet (Ours) | **91.5** | **93.9** | 1.1 M | 4.1 h | 3 ms |

of both training and testing times. Specifically, the training time is only one-third of PointMLP (8.6 h vs. 23 h), and the testing time is also twice as fast as PointMLP (4 ms vs. 7 ms). Compared to the non-parametric Point-NN, our method demonstrates stronger segmentation capabilities, with an Inst. mIoU metric surpassing 14.9%, and the testing time is only one-fifth of Point-NN (4 ms vs. 19 ms).

Furthermore, we also present visualizations of the ground truth and predicted segmentation in Figure 4. Intuitively, the predictions from our PointeNet closely match the ground truth, providing qualitative evidence for the validity of our method.

Table 3: **Shape classification on the real-world ScanObjectNN.** We report the class-average accuracy (mAcc) and overall accuracy (OA) on the most challenging variant of the testing set. ∗ means the reproduced results.

| Method | mAcc (%) | OA (%) | Param. | Train Time | Test Time |
|---|---|---|---|---|---|
| PointNet Qi et al. (2017a) | 63.4 | 68.2 | 3.5 M | - | - |
| SpiderCNN Xu et al. (2018) | 69.8 | 73.7 | - | - | - |
| PointNet++ Qi et al. (2017b) | 75.4 | 77.9 | 1.7 M | - | - |
| PointCNN Li et al. (2018) | 75.1 | 78.5 | 15.2 M | - | - |
| GBNet Qiu et al. (2022) | 77.8 | 80.5 | 8.4 M | - | - |
| SimpleView Goyal et al. (2021) | - | 80.5 | - | - | - |
| PointMLP∗ Ma et al. (2022) | 83.2 | 85.0 | 13.2 M | 7.6 h | 4 ms |
| Point-NN∗ Zhang et al. (2023) | - | 64.8 | 0 M | 0 h | 4 ms |
| PointeNet (Ours) | **86.4** | **87.9** | 4.7 M | 12.3 h | 5 ms |

*4.3. Shape Classification on ModelNet40*

We also evaluate PointeNet and compare our method with several recent works, including PointNet Qi et al. (2017a), PointNet++ Qi et al. (2017b), PointCNN Li et al. (2018), PointConv Wu et al. (2019), KPConv Thomas et al. (2019), PointTransformer Zhao et al. (2021), GBNet Qiu et al. (2022), PointMLP Ma et al. (2022) and Point-NN Zhang et al. (2023), for the 3D shape classification task on the synthetic ModelNet40 Wu et al. (2015) benchmark. Following the standard practice in the community, we report the class-average accuracy (mAcc) and overall accuracy (OA) on the testing set in Table 2.

Our PointeNet outperforms PointMLP by 0.6% in OA (93.9% vs. 93.3%) and 0.7% in mAcc (91.5% vs. 90.8%), all while requiring 12 times fewer parameters (13.2 M vs. 1.1 M). Moreover, in terms of training and testing times, our method is more cost-effective, with half the training time (4.1 h vs. 11 h) and approximately twice the testing speed (3 ms vs. 5 ms) compared to PointMLP. Compared to Point-NN, our method not only demonstrates superior performance in OA (93.9% vs. 81.3%) but also achieves faster

testing speeds (3 ms vs. 4 ms). These experimental results validate the efficiency and superiority of our method in shape classification tasks.

### 4.4. Shape Classification on ScanObjectNN

Although our method has already demonstrated its efficiency in shape classification tasks on the synthetic ModelNet40 dataset, to thoroughly validate this capability, we extend our evaluation to the ScanObjectNN benchmark Uy et al. (2019), where we compare PointeNet with several recent methods, including PointNet Qi et al. (2017a), PointNet++ Qi et al. (2017b), PointCNN Li et al. (2018), SpiderCNN Xu et al. (2018), GBNet Qiu et al. (2022), SimpleView Goyal et al. (2021), PointMLP Ma et al. (2022) and Point-NN Zhang et al. (2023).

As shown in Table 3, our PointMLP surpasses all methods with a significant improvement in both class mean accuracy (mAcc) and overall accuracy (OA). Notably, despite utilizing only one-third of the parameters (4.7 M vs. 13.2 M) employed by PointMLP, our method achieves significantly superior performance. Specifically, we surpass PointMLP by more than 3.2% in mAcc (86.4% vs. 83.2%) and 2.9% in OA (87.9% vs. 85.0%). Compared to Point-NN, our method also maintains a significant advantage. In addition, it's worth noting that our method achieves a smaller gap between class average accuracy (mAcc) and overall accuracy (OA) compared to PointMLP (1.5% vs. 1.8%). This observation suggests that PointeNet avoids bias towards any specific category, demonstrating commendable robustness. Through these results on the real-world ScanObjectNN, the efficiency and superiority of our method in the shape classification task have been thoroughly validated.

Table 4: Results of PointRCNN with and without PointeNet on the KITTI val set. The best performance value is in **bold**, second-best is underlined. ∗ means the reproduced results.

| Method | Car (3D $AP_{R40}$) | | | Pedestrian (3D $AP_{R40}$) | | | Cyclist (3D $AP_{R40}$) | | |
|---|---|---|---|---|---|---|---|---|---|
| | Easy | Mod. | Hard | Easy | Mod. | Hard | Easy | Mod. | Hard |
| SECOND Yan et al. (2018) | 90.55 | 81.61 | 78.56 | 55.94 | 51.15 | 46.17 | 82.97 | 66.74 | 62.78 |
| PointPillars Lang et al. (2019) | 87.75 | 78.41 | 75.19 | 57.30 | 51.42 | 46.87 | 81.57 | 62.93 | 58.98 |
| PV-RCNN Shi et al. (2020) | 92.10 | **84.36** | **82.48** | 64.26 | 56.67 | 51.91 | 88.88 | 71.95 | 66.78 |
| PC-RGNN Zhang et al. (2021) | 90.94 | 81.43 | 80.45 | - | - | - | - | - | - |
| Voxel-RCNN Deng et al. (2021) | 91.72 | 83.19 | 78.60 | - | - | - | - | - | - |
| PointRCNN∗ Shi et al. (2019) | 91.28 | 80.78 | 78.37 | 65.12 | 56.43 | 49.24 | 90.02 | **72.42** | 67.44 |
| PointRCNN + PointeNet | **92.34** | 83.13 | 80.78 | **66.00** | **58.62** | **52.09** | **92.99** | **73.97** | **69.50** |
| *Improvement* | **1.06** | **2.35** | **2.41** | **0.88** | **2.19** | **2.85** | **2.97** | **1.55** | **2.06** |

Table 5: Ablation results of the proposed MLGA on the ShapeNetPart testing set. We analyze the effect of the multivariate geometries, including spatial neighboring, normal, and curvature, as well as the multivariate feature aggregation strategy (i.e., MAA or concatenation), on PointeNet. Note that the embedding dimension of PointeNet is set to 36.

| Geometries | | | Aggregation | | Cls. | Inst. | Param. | Train | Test |
|---|---|---|---|---|---|---|---|---|---|
| Distrib. | Normal | Curvature | MAA | Concat. | mIoU (%) | mIoU (%) | | Time | Time |
| ✓ | | | ✓ | | 82.98 | 94.76 | 4.6 M | 5.8 h | 3 ms |
| ✓ | ✓ | | ✓ | | 83.47 | 84.94 | 5.5 M | 6.1 h | 3 ms |
| ✓ | | ✓ | ✓ | | 83.48 | 85.13 | 5.5 M | 6.1 h | 3 ms |
| ✓ | ✓ | ✓ | | ✓ | 83.65 | 85.09 | 4.6 M | 6.1 h | 3 ms |
| ✓ | ✓ | ✓ | ✓ | | **83.73** | **85.17** | 5.5 M | 6.2 h | 3 ms |

### 4.4.1. 3D Oject Detection on KITTI

To validate the practicality and efficiency of our method in real-world, scene-level applications, we utilize the classic PointRCNN as a baseline and evaluate it on the real-world autonomous driving dataset KITTI Geiger et al. (2012). We compare PointRCNN with several established 3D object detection networks,

Table 6: Ablation results of the proposed DSE on the KITTI val set. We analyze the effect of MGE and DSE (with/without the distance factor $d$) on PointeNet.

| PointRCNN | MGE | DSE | | Car (3D AP$_{R40}$) | | |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| | | w/ $d$ | w/o $d$ | Easy | Mod. | Hard |
| ✓ | | | | 91.28 | 80.78 | 78.37 |
| ✓ | ✓ | | | 91.98 | 82.35 | 80.39 |
| ✓ | ✓ | ✓ | | 92.04 | 82.51 | 80.49 |
| ✓ | ✓ | ✓ | ✓ | **92.34** | **83.13** | **80.78** |

including SECOND Yan et al. (2018), PointPillars Lang et al. (2019), PV-RCNN Shi et al. (2020), PC-RGNN Zhang et al. (2021), and Voxel-RCNN Deng et al. (2021).

From the results of Table 4, we can observe that our PointFPN significantly improves the performance of PointRCNN in all three categories: *Car*, *Pedestrian*, and *Cyclist*, with the highest improvement approaching 3%. Moreover, PointRCNN with the addition of PointeNet also has superior performance compared to the PV-RCNN, which combines the advantageous features of Voxel and Point. From these results, it is evident that our method performs exceptionally well in real-world application scenarios and exhibits high efficiency.

### 4.5. Ablation Study

#### 4.5.1. Effects of MLGA

Our multivariate geometric encoding (MGE) module consists of non-parametric FPS, k-NN, and a multivariate local geometric aggregation (MLGA) module with a minimal number of learnable parameters, where MLGA is crucial. The proposed MLGA captures multivariate geometric features, including spatial neighboring, curvature, and normal information, and adaptively aggregates these geometric features. To thoroughly validate the effectiveness of these geometric features and the proposed multivariate adaptive aggregation (MAA) module, we conduct comprehensive experiments, and the results are presented in Table 5. From the results in the first three rows of Table 2, it can be observed that spatial geometry, curvature, and normal information, these three geometric features, all make significant contributions to PointeNet. The comparative results in the fourth and fifth rows also demonstrate that these geometries, when adaptively aggregated through our MAA module, are more effective compared to directly concatenating them together.

#### 4.5.2. Effects of DSE

We conduct ablation experiments on the challenging *Car* category of the KITTI dataset to analyze the effectiveness of MGE and DSE (with/without the distance factor $d$) in PointeNet. From the results of Table 6, the multivariate geometric features output by MGE significantly improve the performance of PointRCNN, especially with a notable enhancement of 2.02% on the easy level. Without the distance factor to adaptively adjust the segmentation difficulty for semantic features, there is a slight improvement of 0.1% across the three difficulty levels. When DSE is improved with the distance factor, PointeNet assists PointRCNN in further improvement, particularly enhancing performance by 0.62% on the moderate level.

## 5. Conclusion

While prior methods in point cloud analysis have demonstrated promising success, they often face challenges associated with increased complexity, limiting their applicability to real-world scene-level scenarios and posing efficiency concerns. In response to these challenges, this paper introduces a lightweight network specifically optimized for scalability in real-world autonomous driving applications, offering a more efficient approach to point cloud analysis. Our proposed network leverages non-learnable components, such as FPS and K-NN, and introduces a minimal number of learnable parameters through the proposed Multivariate Geometric Encoding (MGE) module. This module adeptly captures and adaptively aggregates multivariate geometric features, ensuring a lightweight yet comprehensive representation of the point cloud. Furthermore, we incorporate a Distance-Aware Semantic Enhancement (DSE) module tailored for autonomous driving

scenarios, capturing distance-aware semantic features. The integration of both multivariate geometric and semantic features contributes to a more efficient and superior performance in point cloud analysis. In future work, we aim to explore more potential point cloud features to further enhance the capabilities of our lightweight network.

## References

Deng, J., Shi, S., Li, P., Zhou, W., Zhang, Y., Li, H., 2021. Voxel r-cnn: Towards high performance voxel-based 3d object detection, in: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 1201–1209.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N., 2021. An image is worth 16x16 words: Transformers for image recognition at scale, in: 9th International Conference on Learning Representations.

Felzenszwalb, P.F., Huttenlocher, D.P., 2004. Efficient graph-based image segmentation. Int. J. Comput. Vis. 59, 167–181.

Geiger, A., Lenz, P., Urtasun, R., 2012. Are we ready for autonomous driving? the KITTI vision benchmark suite, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 3354–3361.

Goyal, A., Law, H., Liu, B., Newell, A., Deng, J., 2021. Revisiting point cloud shape classification with a simple and effective baseline, in: Proceedings of the 38th International Conference on Machine Learning, pp. 3809–3820.

Guo, M., Cai, J., Liu, Z., Mu, T., Martin, R.R., Hu, S., 2021. PCT: point cloud transformer. Comput. Vis. Media 7, 187–199.

He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778.

Klokov, R., Lempitsky, V.S., 2017. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models, in: IEEE International Conference on Computer Vision, pp. 863–872.

Lang, A.H., Vora, S., Caesar, H., Zhou, L., Yang, J., Beijbom, O., 2019. Pointpillars: Fast encoders for object detection from point clouds, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 12697–12705.

Li, Y., Bu, R., Sun, M., Wu, W., Di, X., Chen, B., 2018. Pointcnn: Convolution on x-transformed points, in: Advances in Neural Information Processing Systems 31, pp. 828–838.

Liang, T., Xie, H., Yu, K., Xia, Z., Lin, Z., Wang, Y., Tang, T., Wang, B., Tang, Z., 2022. Bevfusion: A simple and robust lidar-camera fusion framework, in: Advances in Neural Information Processing Systems, pp. 10421–10434.

Lin, T., Goyal, P., Girshick, R.B., He, K., Dollár, P., 2020. Focal loss for dense object detection. IEEE Trans. Pattern Anal. Mach. Intell. 42, 318–327.

Lin, Z., Huang, S., Wang, Y.F., 2022. Learning of 3d graph convolution networks for point cloud analysis. IEEE Trans. Pattern Anal. Mach. Intell. 44, 4212–4224.

Liu, Y., Fan, B., Xiang, S., Pan, C., 2019. Relation-shape convolutional neural network for point cloud analysis, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 8895–8904.

Ma, X., Qin, C., You, H., Ran, H., Fu, Y., 2022. Rethinking network design and local geometry in point cloud: A simple residual MLP framework, in: The Tenth International Conference on Learning Representations.

Qi, C.R., Litany, O., He, K., Guibas, L.J., 2019. Deep hough voting for 3d object detection in point clouds, in: IEEE International Conference on Computer Vision, pp. 9276–9285.

Qi, C.R., Su, H., Mo, K., Guibas, L.J., 2017a. Pointnet: Deep learning on point sets for 3d classification and segmentation, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 77–85.

Qi, C.R., Yi, L., Su, H., Guibas, L.J., 2017b. Pointnet++: Deep hierarchical feature learning on point sets in a metric space, in: Advances in Neural Information Processing Systems, pp. 5099–5108.

Qiu, S., Anwar, S., Barnes, N., 2022. Geometric back-projection network for point cloud classification. IEEE Trans. Multim. 24, 1943–1955.

Ran, H., Zhuo, W., Liu, J., Lu, L., 2021. Learning inner-group relations on point clouds, in: 2021 IEEE/CVF International Conference on Computer Vision, pp. 15457–15467.

Shi, S., Guo, C., Jiang, L., Wang, Z., Shi, J., Wang, X., Li, H., 2020. PV-RCNN: point-voxel feature set abstraction for 3d object detection, in: 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2020, Seattle, WA, USA, June 13-19, 2020, pp. 10526–10535.

Shi, S., Wang, X., Li, H., 2019. Pointrcnn: 3d object proposal generation and detection from point cloud, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–779.

Su, H., Jampani, V., Sun, D., Maji, S., Kalogerakis, E., Yang, M., Kautz, J., 2018. Splatnet: Sparse lattice networks for point cloud processing, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, pp. 2530–2539.

Thomas, H., Qi, C.R., Deschaud, J., Marcotegui, B., Goulette, F., Guibas, L.J., 2019. Kpconv: Flexible and deformable convolution for point clouds, in: 2019 IEEE/CVF International Conference on Computer Vision, pp. 6410–6419.

Uy, M.A., Pham, Q., Hua, B., Nguyen, D.T., Yeung, S., 2019. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data, in: 2019 IEEE/CVF International Conference on Computer Vision, pp. 1588–1597.

Wang, X., Girshick, R.B., Gupta, A., He, K., 2018. Non-local neural networks, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, pp. 7794–7803.

Wang, Y., Sun, Y., Liu, Z., Sarma, S.E., Bronstein, M.M., Solomon, J.M., 2019. Dynamic graph CNN for learning on point clouds. ACM Trans. Graph. 38, 146:1–146:12.

Wu, W., Qi, Z., Li, F., 2019. Pointconv: Deep convolutional networks on 3d point clouds, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019, pp. 9621–9630.

Wu, Z., Song, S., Khosla, A., Yu, F., Zhang, L., Tang, X., Xiao, J., 2015. 3d shapenets: A deep representation for volumetric shapes, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 1912–1920.

Xu, M., Ding, R., Zhao, H., Qi, X., 2021. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds, in: IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021, pp. 3173–3182.

Xu, Y., Fan, T., Xu, M., Zeng, L., Qiao, Y., 2018. Spidercnn: Deep learning on point sets with parameterized convolutional filters, in: Computer Vision - ECCV 2018 - 15th European Conference, pp. 90–105.

Yan, Y., Mao, Y., Li, B., 2018. Second: Sparsely embedded convolutional detection. Sensors 18, 3337.

Yang, Z., Sun, Y., Liu, S., Jia, J., 2020. 3dssd: Point-based 3d single stage object detector, in: IEEE Conference on Computer Vision and Pattern Recognition, pp. 11037–11045.

Yang, Z., Sun, Y., Liu, S., Shen, X., Jia, J., 2019. STD: sparse-to-dense 3d object detector for point cloud, in: 2019 IEEE/CVF International Conference on Computer Vision, pp. 1951–1960.

Yi, L., Kim, V.G., Ceylan, D., Shen, I., Yan, M., Su, H., Lu, C., Huang, Q., Sheffer, A., Guibas, L.J., 2016. A scalable active framework for region annotation in 3d shape collections. ACM Trans. Graph. 35, 210:1–210:12.

Zhang, R., Wang, L., Wang, Y., Gao, P., Li, H., Shi, J., 2023. Parameter is not all you need: Starting from non-parametric networks for 3d point cloud analysis. arXiv preprint arXiv:2303.08134 .

Zhang, W., Zhou, H., Sun, S., Wang, Z., Shi, J., Loy, C.C., 2019. Robust multi-modality multi-object tracking, in: 2019 IEEE/CVF International Conference on Computer Vision, pp. 2365–2374.

Zhang, Y., Huang, D., Wang, Y., 2021. Pc-rgnn: Point cloud completion and graph neural network for 3d object detection, in: Proceedings of the AAAI Conference on artificial intelligence, pp. 3430–3437.

Zhao, H., Jiang, L., Jia, J., Torr, P.H.S., Koltun, V., 2021. Point transformer, in: 2021 IEEE/CVF International Conference on Computer Vision, pp. 16239–16248.