

Parameter is Not All You Need: Starting from Non-Parametric Networks for 3D Point Cloud Analysis

Renrui Zhang^{1,3}, Liuhui Wang^{2,4}, Yali Wang¹, Peng Gao¹, Hongsheng Li³, Jianbo Shi^{4,5}

¹Shanghai Artificial Intelligence Laboratory ²Peking University

³The Chinese University of Hong Kong ⁴Heisenberg Robotics ⁵University of Pennsylvania

{zhangrenrui, gaopeng}@pjlab.org.cn, jshi@seas.upenn.edu
wangliuhui0401@pku.edu.cn, hsli@ee.cuhk.edu.hk

Abstract

We present a Non-parametric Network for 3D point cloud analysis, **Point-NN**, which consists of purely non-learnable components: farthest point sampling (FPS), k -nearest neighbors (k -NN), and pooling operations, with trigonometric functions. Surprisingly, it performs well on various 3D tasks, requiring no parameters or training, and even surpasses existing fully trained models. Starting from this basic non-parametric model, we propose two extensions. First, Point-NN can serve as a base architectural framework to construct Parametric Networks by simply inserting linear layers on top. Given the superior non-parametric foundation, the derived **Point-PN** exhibits a high performance-efficiency trade-off with only a few learnable parameters. Second, Point-NN can be regarded as a plug-and-play module for the already trained 3D models during inference. Point-NN captures the complementary geometric knowledge and enhances existing methods for different 3D benchmarks without re-training. We hope our work may cast a light on the community for understanding 3D point clouds with non-parametric methods. Code is available at <https://github.com/ZrrSkywalker/Point-NN>.

1. Introduction

Point cloud 3D data processing is a foundational operation in autonomous driving [4, 28, 44], scene understanding [1, 3, 64, 89], and robotics [7, 43, 53]. Point clouds contain unordered points discretely depicting object surfaces in 3D space. Unlike grid-based 2D images, they are distribution-irregular and permutation-invariant, which leads to non-trivial challenges for algorithm designs.

Since PointNet++ [46], the prevailing trend has been adding advanced local operators and scaled-up learnable pa-

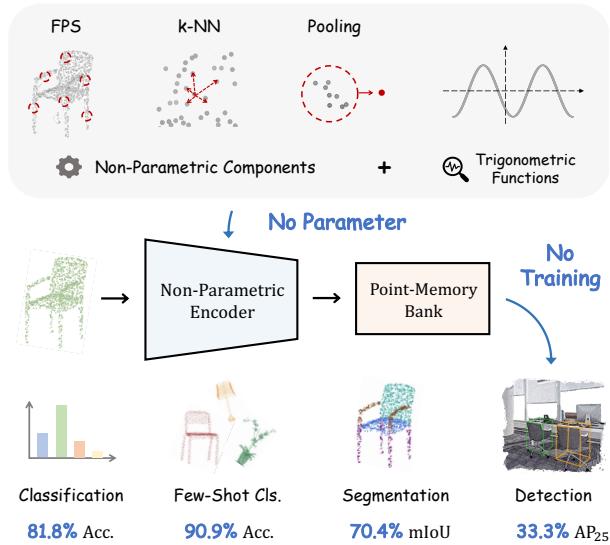


Figure 1. **The Pipeline of Non-Parametric Networks.** Point-NN is constructed by the basic non-parametric components without any learnable operators. Free from training, Point-NN can achieve favorable performance on various 3D tasks.

rameters. Instead of max pooling for feature aggregation, mechanisms are proposed to extract local geometries, e.g., adaptive point convolutions [35, 38, 61, 70, 71] and graph-like message passing [19, 65, 87]. The performance gain also rises from scaling up the number of parameters, e.g., KPConv [61]'s 14.3M and PointMLP [39]'s 12.6M, is much larger than PointNet++'s 1.7M. This trend has increased network complexity and computational resources.

Instead, the non-parametric framework underlying all the learnable modules remains nearly the same since PointNet++, including farthest point sampling (FPS), k -Nearest Neighbors (k -NN), and pooling operations. Given that few works have investigated their efficacy, we ask the question:

can we achieve high 3D point cloud analysis performance using only these non-parametric components?

We present a Non-parametric Network, termed Point-NN, which is constructed by the aforementioned non-learnable components. Point-NN, as shown in Figure 1, consists of a non-parametric encoder for 3D feature extraction and a point-memory bank for task-specific recognition. The multi-stage encoder applies FPS, k -NN, trigonometric functions, and pooling operations to progressively aggregate local geometries, producing a high-dimensional global vector for the point cloud. We only adopt simple trigonometric functions to reveal local spatial patterns at every pooling stage without learnable operators. Then, we adopt the non-parametric encoder of Point-NN to extract the training-set features and cache them as the point-memory bank. For a test point cloud, the bank outputs task-specific predictions via naive feature similarity matching, which validates the encoder’s discrimination ability.

Free from any parameters or training, Point-NN unexpectedly achieves favorable performance on various 3D tasks, e.g., shape classification, part segmentation and 3D object detection, indicating the strength of the long-ignored non-parametric operations. Compared to existing parametric methods, Point-NN even surpasses the fully trained 3DmFV [2] by **+2.9%** classification accuracy on OBJ-BG split of ScanObjectNN [62]. Especially for few-shot classification, Point-NN significantly exceeds PointCNN [33] and other models [45, 46] by more than +20% accuracy, indicating its superiority in low-data regimes. *Starting from this simple-but-effective Point-NN, we propose two approaches by leveraging the non-parametric components to benefit 3D point cloud analysis.*

First, Point-NN can serve as an architectural precursor to construct Parametric Networks, termed as Point-PN, shown in Figure 2 (a). As we have fully optimized the non-parametric framework, the Point-PN can be simply derived by inserting linear layers into every stage of Point-NN. Point-PN contains no complicated local operators, but only linear layers and trigonometric functions inherited from Point-NN. Experiments show that Point-PN can achieve competitive performance with a small number of parameters, e.g., 87.1% classification accuracy with 0.8M on the hardest split of ScanObjectNN [62].

Second, Point-NN can be regarded as a plug-and-play module to enhance the already trained 3D models without re-training, as shown in Figure 2 (b). We directly fuse the predictions between Point-NN and off-the-shelf 3D models during inference by linear interpolation. Given the training-free manner, Point-NN mainly focuses on low-level 3D structural signals by trigonometric functions, which provides complementary knowledge to the high-level semantics of existing 3D models. On different tasks, Point-NN exhibits consistent performance boost, e.g., +2.0% classifi-

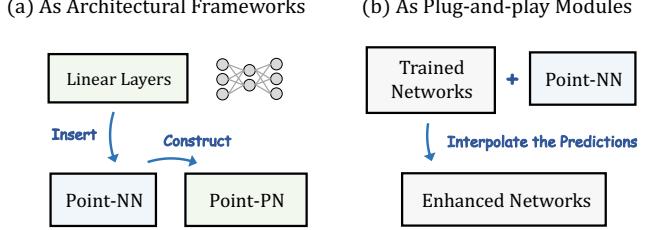


Figure 2. **Two Applications of Point-NN.** (a) As a non-parametric framework to construct parametric networks by simply inserting linear layers. (b) As a plug-and-play module to enhance already trained networks without re-training.

cation accuracy on ScanObjectNN [62] and +11.02% detection AR₂₅ on ScanNetV2 [9].

Our contributions are summarized below:

- We propose to revisit the non-learnable components in 3D networks, and, *for the first time*, develop a non-parametric method, Point-NN, for 3D analysis.
- By using Point-NN as a basic framework, we introduce its parameter-efficient derivative, Point-PN, which exerts superior performance without advanced operators.
- As a plug-and-play module, the complementary Point-NN can boost off-the-shelf trained 3D models for various 3D tasks directly during inference.

2. Non-Parametric Networks

In this section, we first investigate the basic non-parametric components in existing 3D models (Sec. 2.1). Then, we integrate them into Point-NN, which consists of a non-parametric encoder (Sec. 2.2) and a point-memory bank (Sec. 2.3). Finally, we introduce how to apply Point-NN for various 3D tasks (Sec. 2.4).

2.1. Basic Components

Underlying most of the point cloud processing networks [39, 46] are a series of non-parametric components, i.e., farthest point sampling (FPS), k -Nearest Neighbors (k -NN), and pooling operations. These building blocks are non-learnable during training and iteratively stacked into multiple stages to construct a pyramid hierarchy.

For each stage, we denote the input point cloud representation from the last stage as $\{p_i, f_i\}_{i=1}^M$, where $p_i \in \mathbb{R}^{1 \times 3}$ and $f_i \in \mathbb{R}^{1 \times C}$ denote the coordinate and feature of point i . First, FPS is adopted to downsample the point number from M to $\frac{M}{2}$ by selecting a subset of local center points,

$$\{p_c, f_c\}_{c=1}^{\frac{M}{2}} = \text{FPS} \left(\{p_i, f_i\}_{i=1}^M \right). \quad (1)$$

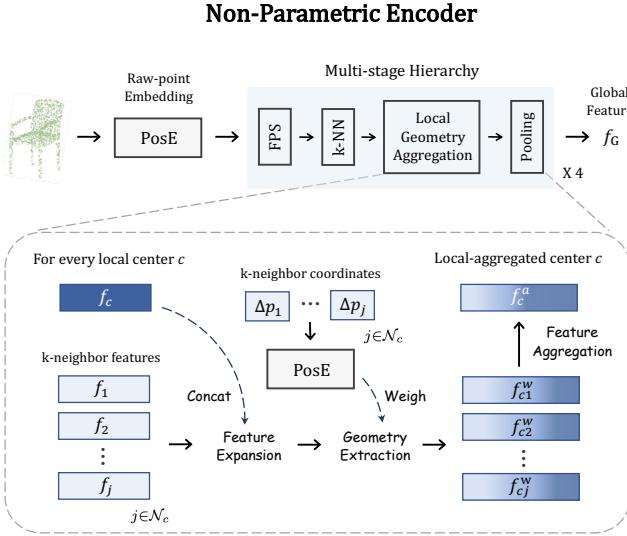


Figure 3. **Non-Parametric Encoder of Point-NN.** We first utilize trigonometric functions, denoted as $\text{PosE}(\cdot)$, to encode raw points into high-dimensional vectors, and then adopt non-learnable operations to hierarchically aggregate local features.

Then, k -NN, or ball query, is responsible for grouping k spatial neighbors for each center c from the original M points, which forms the local 3D region,

$$\mathcal{N}_c = k\text{-NN}\left(p_c, \{p_i\}_{i=1}^M\right), \quad (2)$$

where $\mathcal{N}_c \in \mathbb{R}^{k \times 1}$ denotes the indices for k nearest neighbors. On top of this, the geometric patterns for each local neighborhood \mathcal{N}_c are extracted by the delicate learnable operator, $\Phi(\cdot)$, and finally aggregated by max pooling, $\text{MaxP}(\{\cdot\})$. We formulate it as

$$f_c^a = \text{MaxP}\left(\{\Phi(f_c, f_j)\}_{j \in \mathcal{N}_c}\right). \quad (3)$$

The derived $\{p_c, f_c^a\}_{i=1}^{\frac{M}{2}}$ is then fed into the next network stage to progressively capture 3D geometries with enlarging receptive fields.

To verify the non-parametric efficacy independently from the learnable $\Phi(\cdot)$, we present Point-NN, a network purely constructed by these non-learnable basic components along with simple trigonometric functions for 3D coordinate encoding. Point-NN is composed of a non-parametric encoder $\text{NPEnc}(\cdot)$ and a point-memory bank $\text{PoM}(\cdot)$. Given an input point cloud P for shape classification, the encoder summarizes its high-dimensional global feature f_G , and the bank produces the classification logits by similarity matching. We formulate it as

$$f_G = \text{EncNP}(P); \text{ logits} = \text{PoM}(f_G), \quad (4)$$

where $f_G \in \mathbb{R}^{1 \times C_G}$ and the $\text{logits} \in \mathbb{R}^{1 \times K}$ denote the predicted possibility for K categories in the dataset.

2.2. Non-Parametric Encoder

As shown in Figure 3, the non-parametric encoder conducts initial embedding to transform the raw XYZ coordinates of P into high-dimensional vectors, and progressively aggregates local patterns via the multi-stage hierarchy.

Raw-point Embedding. To achieve feature embedding without learnable layers, we refer to the positional encoding in Transformer [63] and extend it for non-parametric 3D encoding. For a raw point i with $p_i = (x_i, y_i, z_i) \in \mathbb{R}^{1 \times 3}$, we utilize trigonometric functions to embed it into a C_I -dimensional vector,

$$\text{PosE}(p_i) = \text{Concat}(f_i^x, f_i^y, f_i^z) \in \mathbb{R}^{1 \times C_I}, \quad (5)$$

where $f_i^x, f_i^y, f_i^z \in \mathbb{R}^{1 \times \frac{C_I}{3}}$ denote the embeddings of three axes, and C_I denotes the initial feature dimension. Taking f_i^x as an example, for the channel index $m \in [0, \frac{C_I}{6}]$:

$$f_i^x[2m] = \sin(\alpha x_i / \beta^{\frac{6m}{C_I}}), \quad (6)$$

$$f_i^x[2m+1] = \cos(\alpha x_i / \beta^{\frac{6m}{C_I}}), \quad (7)$$

where α, β control the magnitude and wavelengths, respectively. By the inherent nature of trigonometric functions, the transformed vectors can well encode relative positional information between different points and capture fine-grained structural variation of 3D shapes.

Local Geometry Aggregation. Based on the embedding, we adopt a 4-stage network architecture to hierarchically aggregate spatial local features. After the ordinal FPS and k -NN illustrated in Section 2.1, we discard any learnable operators $\Phi(\cdot)$, and simply utilize trigonometric functions $\text{PosE}(\cdot)$ to reveal the local patterns. In detail, for each center point $\{p_c, f_c\}$ and its neighborhood $\{p_j, f_j\}_{j \in \mathcal{N}_c}$, we aim to achieve there goals. **(1) Feature Expansion.** As the network stage goes deeper, each point feature is assigned with larger receptive field and requires higher feature dimension to encode 3D semantics. We conduct such feature expansion by simply concatenating the neighbor feature p_j with the center feature p_c along the feature dimension,

$$f_{cj} = \text{Concat}(f_c, f_j), \quad \text{for } j \in \mathcal{N}_c, \quad (8)$$

where f_{cj} denotes the expanded feature of each neighbor.

(2) Geometry Extraction. To indicate the spatial distribution of k neighbors within the local region, we weigh each f_{cj} by the relative positional encoding. We normalize their coordinates by the mean and standard deviation, denoted as $\{\Delta p_j\}_{j \in \mathcal{N}_c}$, and embed them via Eq. (5). Then, the k -neighbor features are weighted as

$$f_{cj}^w = (f_{cj} + \text{PosE}(\Delta p_j)) \odot \text{PosE}(\Delta p_j), \quad (9)$$

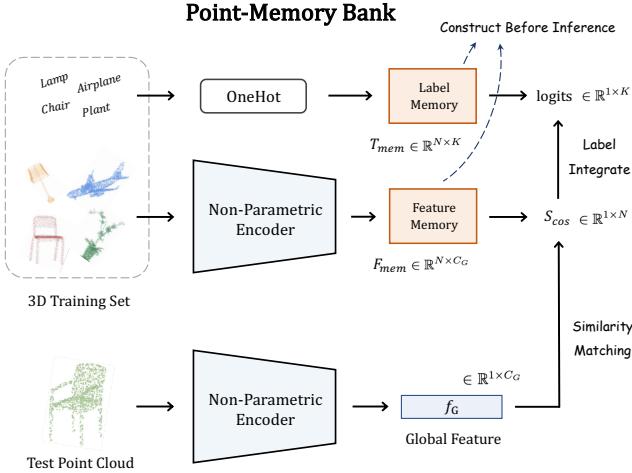


Figure 4. **Point-Memory Bank of Point-NN.** We construct the memory bank by caching training-set features via the non-parametric encoder. Then, the test point cloud is simply classified by similarity matching without training.

where \odot denotes element-wise multiplication. In this way, the local geometry of the region can be implicitly encoded into the features without any learnable parameters. **(3) Feature Aggregation.** After the weighing, we utilize both max and average pooling for local feature aggregation,

$$f_c^a = \text{MaxP} \left(\{f_{c,j}^w\}_{j \in \mathcal{N}_c} \right) + \text{AveP} \left(\{f_{c,j}^w\}_{j \in \mathcal{N}_c} \right), \quad (10)$$

where $\text{MaxP}(\{\cdot\})$, $\text{AveP}(\{\cdot\})$ are permutation-invariant and summarize neighboring features from different aspects. Here, we obtain the local-aggregated centers $\{f_c^a, p_c\}_{i=1}^{\frac{M}{2}}$, which would be fed into the next stage of Point-NN. Finally, after all 4 stages, we apply the two pooling operations to integrate the features and acquire a global representation f_G with C_G feature dimension of the input point cloud.

2.3. Point-Memory Bank

Instead of using the traditional learnable classification head, our Point-NN adopts a point-memory bank to involve sufficient category knowledge from the 3D training set. As shown in Figure 4, the bank is first constructed by the non-parametric encoder in a training-free manner, and then outputs the prediction by similarity matching during inference.

Memory Construction. The point memory consists of a feature memory F_{mem} and a label memory T_{mem} . Taking the task of shape classification as an example, we suppose the given training set contains N point clouds, $\{P_n\}_{n=1}^N$, of K categories. Via the aforementioned non-parametric encoder, we encode all N global features and convert their ground-truth labels $\{t_n\}_{n=1}^N$ as one-hot encoding. We then

cache them as two matrices by concatenating along the inter-sample dimension as

$$F_{mem} = \text{Concat} \left(\{\text{EncNP}(P_n)\}_{n=1}^N \right), \quad (11)$$

$$T_{mem} = \text{Concat} \left(\{\text{OneHot}(t_n)\}_{n=1}^N \right), \quad (12)$$

where $F_{mem} \in \mathbb{R}^{N \times C_G}$ and $T_{mem} \in \mathbb{R}^{N \times K}$. Tagged by T_{mem} , the memory F_{mem} can be regarded as the encoded category knowledge for the 3D training set. Features tagged with the same label unitedly describe the characteristics of the same category, and the inter-class discrimination can also be reflected by the embedding-space distances.

Similarity-based Prediction. For a test point cloud, we also utilize the non-parametric encoder to extract its global feature as $f_G^t \in \mathbb{R}^{1 \times C_G}$, which is within the same embedding space as the feature memory F_{mem} . Then, the classification is simply accomplished by two matrix multiplications via the constructed bank. Firstly, we calculate the cosine similarity between the test feature and F_{mem} by

$$S_{cos} = \frac{f_G^t F_{mem}^T}{\|f_G^t\| \cdot \|F_{mem}\|} \in \mathbb{R}^{1 \times N}, \quad (13)$$

where S_{cos} denotes the semantic correlation of the test point cloud and every N training samples. On top of that, we integrate the one-hot labels in the label memory T_{mem} weighted by S_{cos} , formulated as

$$\text{logits} = \varphi(S_{cos} T_{mem}) \in \mathbb{R}^{1 \times K}, \quad (14)$$

where $\varphi(x) = \exp(-\gamma(1-x))$ serves as an activation function from [77]. In S_{cos} , the more similar feature memory of higher score contributes more to the final classification logits, and vice versa. By such similarity-based label integration, the point-memory bank can adaptively discriminate different point cloud instances without any training.

2.4. Other 3D Tasks

Except for shape classification, our Point-NN can also be extended for part segmentation and 3D object detection with tasks-specific modifications.

Part Segmentation. Other than extracting global features, the task of part segmentation requires to classify each input point. Therefore, we append a symmetric non-parametric decoder after the encoder, which progressively upsamples the encoded point features into the input point number. In every stage of the decoder, we reversely propagate the features of center points to their k neighbors in a non-parametric manner. For the point-memory bank, we first apply the non-parametric encoder and decoder to extract all point-wise features of the training set. To save the

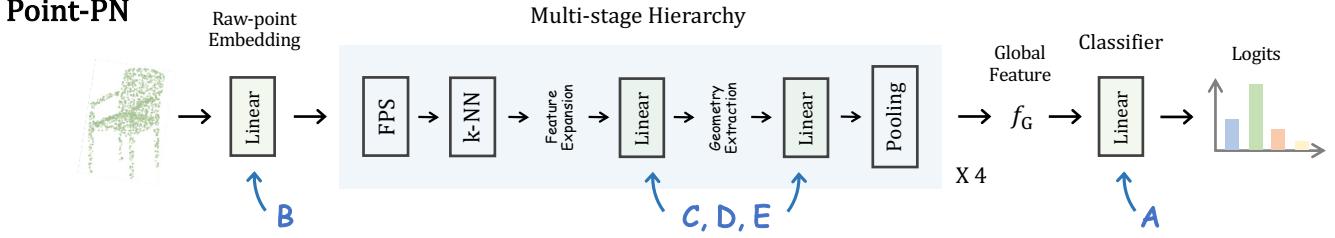


Figure 5. **The Pipeline of Point-PN.** Given the non-parametric framework of Point-NN, we simply construct the parametric derivative, Point-PN, by inserting linear layers into every stage of the encoder. Performance gain of using linear layers of A~E is shown in Table 1.



Figure 6. **Complementary Characteristics of Point-PN.** We visualize the feature responses for the trained PointNet++ [46] and Point-PN, where darker colors indicate higher responses. As shown, they focus on different spatial structures with complementary 3D patterns.

GPU memory, we average the features of points with the same part label in an object, and only cache such aggregated part-wise features with part labels as F_{mem}, T_{mem} .

3D Object Detection. Given category-agnostic 3D proposals from a pre-trained 3D region proposal network [11, 42], Point-PN can be utilized as a non-parametric classification head for object detection. Similar to shape classification, we also adopt a pooling operation after the encoder to obtain global features of the detected objects. Differently, we sample the point cloud within each ground-truth 3D bounding box in the training set, and encode the object-wise features as the feature memory F_{mem} . Specifically, we do not normalize the point coordinates for each object as the pre-processing like other 3D tasks, which is to preserve the 3D positional information of objects in the original scene.

3. Starting from Point-NN

In this section, we introduce two promising applications for Point-PN, which fully unleash the potentials of non-parametric components for 3D point cloud analysis.

3.1. As Architectural Frameworks

Point-NN can be extended into learnable parametric networks (Point-PN) without adding complicated operators or too many parameters. As shown in Figure 5 and Table 1, on

Method	Raw Embed.	Linear Layers	Classifier	Acc. (%)	Param.
Point-NN	N	-	N	81.8	0.0 M
A	N	-	P	90.3	0.3 M
B	P	-	P	90.8	0.3 M
C	P	0+1	P	93.4	0.5 M
D	P	1+1	P	93.2	0.8 M
E	P	2+2	P	92.9	0.7 M
Point-PN	P	1+2	P	93.8	0.8 M

Table 1. **Step-by-step Construction of Point-PN** on ModelNet40 [69]. ‘N’ or ‘P’ denotes the non-parametric modules or parametric linear layers. ‘Linear Layers’ denotes the number of linear layers inserted into each network stage. Note that we adopt bottleneck layers with a ratio 0.5 for ‘2’ in ‘Linear Layers’

top of Point-PN, we first replace the point-memory bank with a conventional learnable classifier. This lightweight version achieves 90.3% classification accuracy on ModelNet40 [69] with only 0.3M parameters (A). Then, we upgrade the raw-point embedding into parametric linear layers, which improves the performance to 90.8% (B). To better extract multi-scale hierarchy, we append simple linear layers into every stage of the encoder. For each stage, two learnable linear layers are inserted right before or after the *Geometry Extraction* step for capturing higher-level spatial

Method	Split 1	Split 2	Split 3	Param.
3DmFV [2]	68.2	73.8	63.0	-
PointNet [45]	73.3	79.2	68.2	3.5 M
SpiderCNN [73]	77.1	79.5	73.7	-
PointNet++ [46]	82.3	84.3	77.9	1.7 M
DGCNN [65]	82.8	86.2	78.1	1.8 M
PointCNN [33]	86.1	85.5	78.5	-
DRNet [48]	-	-	80.3	-
GBNet [49]	-	-	80.5	8.4 M
SimpleView [18]	-	-	80.5	-
PointMLP [39]	-	-	85.2	12.6 M
Point-NN	71.1	74.9	64.9	0.0 M
Point-PN	91.0	90.2	87.1	0.8 M

Table 2. **Shape Classification on the Real-world ScanObjectNN [62].** We report the accuracy (%) on three official splits of ScanObjectNN: OBJ-BG, OBJ-ONLY and PB-T50-RS. Our Point-NN outperforms the fully trained 3DmFV as marked in blue.

patterns (C, D, E). We observe Point-PN attains the competitive 93.8% accuracy with 0.8M parameters. This final version only contains trigonometric functions for geometry extraction and simple linear layers for feature transformation. This demonstrates that, compared to existing advanced operators or scaled-up parameters, we can alternatively start from a non-parametric framework, i.e., Point-NN, to obtain a powerful and efficient 3D model.

3.2. As Plug-and-play Modules

Considering the training-free characteristic, we propose to regard Point-NN as an inference-time enhancement module, which can boost already trained 3D models without extra re-training. For shape classification, we directly fuse the prediction by linear interpolation, namely, adding the classification logits of Point-NN and off-the-shelf models element-wisely. This design produces the ensemble for two types of knowledge: the low-level structural signals from Point-NN, and the high-level semantics from the trained networks. As visualized in Figure 6, the extracted point cloud features by Point-NN produce high response values around the sharp 3D structures, e.g., the airplane’s wingtips, chair’s legs, and lamp poles. In contrast, the trained PointNet++ focuses more on 3D structures with rich semantics, e.g., airplane’s main body, chair’s bottoms, and lampshades.

For part segmentation and 3D object detection, concurrently running an extra Point-NN to enhance existing models would be expensive in both time and memory. Thus, referring to [82], we directly adopt the encoders of already trained models to extract point cloud features, and only apply our point-memory bank on top for plug-and-play. In this way, we can also achieve performance improvement by

Method	Acc. (%)	Param.	Train Time	Test Speed
PointNet [45]	89.2	3.5 M	-	-
PointNet++ [46]	90.7	1.7 M	3.4 h	521
DGCNN [65]	92.9	1.8 M	2.4 h	617
RS-CNN [36]	92.9	1.3 M	-	-
DensePoint [35]	93.2	-	-	-
PCT [19]	93.2	-	-	-
GBNet [49]	93.8	8.4 M	-	189
CurveNet [70]	93.8	2.0 M	6.7 h	25
PointMLP [39]	94.1	12.6 M	14.4 h	189
Point-NN	81.8	0.0 M	0	275
Point-PN	93.8	0.8 M	3.9 h	1176

Table 3. **Shape Classification on Synthetic ModelNet40 [69].** All compared methods take 1,024 points as input. Train Time and Test Speed (samples/second) are tested on one RTX 3090 GPU. We report the accuracy without the voting strategy.

Method	Inst. mIoU	Param.	Train Time	Test Speed
PointNet [45]	83.7	8.3 M	-	-
PointNet++ [46]	85.1	1.8 M	26.5 h	45
PACConv [19]	86.0	-	-	-
PointMLP [39]	86.1	16.8 M	47.1 h	119
CurveNet [70]	86.6	5.5 M	56.9 h	22
Point-NN	70.4	0.0 M	0	51
Point-PN	86.6	3.9 M	29.0 h	131

Table 4. **Part Segmentation on ShapeNetPart [76].** All compared methods take 2,048 points as input, and are evaluated by mean IoU scores (%) across instances.

leveraging the complementary knowledge between similarity matching and traditional learnable classification heads.

4. Experiments

We conduct extensive experiments to evaluate the efficacy of Point-NN (Sec. 4.2), and its two extensions as Point-PN (Sec. 4.3) and plug-and-play modules (Sec. 4.4).

4.1. Dataset

For **shape classification**, we report the performance on two benchmarks: ModelNet40 [69] and ScanObjectNN [62]. The former is the most popular 3D dataset, which contains 9,843/2,468 synthesized point clouds with 40 categories for training/validation. The latter is more challenging and consists of 11,416/2,882 training/validation point clouds with 15 categories captured from the noisy

Method	5-way		10-way	
	10-shot	20-shot	10-shot	20-shot
DGCNN [65]	31.6	40.8	19.9	16.9
FoldingNet [75]	33.4	35.8	18.6	15.4
PointNet++ [46]	38.5	42.4	23.0	18.8
PointNet [45]	52.0	57.8	46.6	35.2
3D-GAN [66]	55.8	65.8	40.3	48.4
PointCNN [33]	65.4	68.6	46.6	50.0
Point-NN	88.8	90.9	79.9	84.9
<i>Improvements</i>	+23.4	+22.3	+33.3	+34.9

Table 5. **Few-shot Classification on ModelNet40** [69]. We calculate the average accuracy (%) over 10 independent runs. The reported results of existing methods are taken from [55].

real-world scenes. Considering the background and data augmentation, ScanObjectNN is officially split into three subsets: OBJ-BG, OBJ-ONLY, and PB-T50-RS. For **few-shot classification**, we evaluate on the few-shot subset of ModelNet40 with four different settings, 5-way 10-shot, 5-way 20-shot, 10-way 10-shot and 10-way 20-shot. For **part segmentation**, we adopt ShapeNetPart [76] dataset with synthesized 3D shapes of 16 shape and 50 part categories. It consists of 14,007/2,874 samples for training/validation. For **3D object detection**, we conduct the experiments on ScanNetV2 [9], which includes 1,201/312 training/validation scenes with axis-aligned bounding boxes for 18 object categories. Note that, we report all results **without the voting strategy**.

4.2. Point-NN

Shape Classification. The results of Point-NN are reported in the green rows of Table 2 and 3, respectively for the two datasets. As shown, Point-NN attains favorable classification accuracy for both real-world and synthetic point clouds, indicating our effectiveness and generalizability. Surprisingly, Point-NN even surpasses the fully trained 3DmFV [2] by +2.9%, +1.1%, +1.9% on the three splits of ScanObjectNN [62]. By this, we demonstrate that, without any parameters or training, the non-parametric network components with simple trigonometric functions can achieve satisfactory 3D point cloud recognition.

Few-shot Classification. As shown in Table 5, compared to existing trained models, our Point-NN exerts leading few-shot performance and exceeds the second-best PointCNN [33] by significant margins, e.g., +34.9% for the 10-way 20-shot setting. Given limited training samples, traditional networks with learnable parameters severely suffer from the over-fitting issue, while our Point-NN inherently tackles it by the training-free manner. This indicates our

Method	AP ₂₅	AP ₅₀	AR ₂₅	AR ₅₀
VoteNet [11]	57.8	33.5	80.9	51.3
Point-NN	4.5	3.3	80.9	51.3
Point-NN w/o nor.	23.1	16.0	90.0	51.3
3DETR-m [42]	64.6	46.4	77.2	59.2
Point-NN	7.7	5.7	77.2	59.2
Point-NN w/o nor.	33.3	24.7	77.4	59.3

Table 6. **3D Object Detection on ScanNetV2** [9]. We report the mean Average Precision (%) and mean Average Recall (%) with 0.25 and 0.5 IoU thresholds. ‘nor’ denotes to normalize the point coordinates for each object proposal.

superior capacity for 3D recognition in low-data regimes.

Part Segmentation. Equipped with a non-parametric decoder illustrated in Sec. 2.4, we extend Point-NN for part segmentation and report the mIoU scores over instances in the green row of Table 4. The 70.4% mIoU reveals that our non-parametric network can also produce well-performed point-wise features, and capture the discriminative characteristics for fine-grained spatial understanding.

3D Object Detection. Regarding Point-NN as a non-parametric classification head, we utilize two popular 3D detectors, VoteNet [11] and 3DETR-m [42], to provide category-agnostic 3D region proposals. In Table 6, we compare the performance of Point-NN with and without the normalization for object-wise point coordinates (Sec. 2.4). As shown, processing the point coordinates without normalization can largely enhance the AP scores of Point-NN, which preserves more positional cues of objects’ 3D locations in the original scene. Also, we observe slight AR score improvement over the region proposal networks, since the classification logits of Point-NN can affect the 3D NMS post-processing to remove the false positive boxes.

Ablation Study. In Table 7, we investigate the designs of Point-NN’s non-parametric encoder. We first explore the non-parametric embedding used as PosE, where the trigonometric functions (‘Our Sin/cos’) from Transformers [63] perform the best. Then, we experiment with different approaches in *Geometry Extraction* to weigh the k -neighbor features, and observe utilizing ‘Add+Multiply’ can fully reveal the local geometry. In Table 9 for Point-NN’s point-memory bank, we verify that the cosine similarity better exerts the discrimination capacity of the encoder among other distance metrics. In addition, compared to traditional machine learning algorithms, our memory bank can benefit from the non-learnable similarity-based label integration, and exhibit superior classification accuracy.

Embedding Function	Acc (%)	Weighted by PosE	Acc (%)
w/o	68.6	w/o	77.8
NeRF's [41]	70.1	Add	78.3
Fourier's [58]	76.9	Multiply	80.4
Our Sin/cos	81.8	Add+Multiply	81.8

Table 7. **Ablation Study of Non-Parametric Encoder.** We ablate the non-parametric functions for raw-point embedding and local weighing, and experiment different weighing methods for local geometry extraction.

Method	ScanObjNN	+NN	ModelNet40	+NN
Point-NN	64.9	-	81.8	-
PointNet	68.2	+2.2	89.7	+0.4
PointNet++	77.9	+1.2	92.6	+0.5
PCT	-	-	93.2	+0.2
PointMLP	85.2	+2.0	94.1	+0.3

Table 8. **Plug-and-play for Shape Classification.** We report the performance gain (%) on the PB-T50-RS split of ScanObjectNN [62] and ModelNet40 [69].

4.3. Point-PN

Shape Classification. As shown in Table 2 and 3, constructed from Point-NN, the derived Point-PN achieves competitive results for both real-world and synthetic 3D recognition. Compared to the large-scale PointMLP [39] with stacked MLPs of 12.6M parameters, Point-PN only contains simple linear layers and surpasses it by +1.9% accuracy on ScanObjectNN [62] with 16× fewer parameters and 6× faster inference speed. With simple trigonometric functions, Point-PN attains comparable performance to CurveNet [70] with complicated curve-based grouping on ModelNet40 [69], while contains 2.5× fewer parameters and 6× faster inference speed. The superior classification accuracy of Point-PN fully demonstrates the significance of a powerful non-parametric framework.

Part Segmentation. For point-wise segmentation task in Table 4, Point-PN also achieves competitive performance, i.e., 86.6% mIoU, among existing methods. Compared to CurveNet, Point-PN with simple local geometry aggregation can save 28h training time and inference 6× faster.

Ablation Study. In Table 1, we present how to step-by-step construct Point-PN from the non-parametric Point-NN. As illustrated in Sec. 4.2, we insert linear layers into the encoder's every stage before and after the *Geometry Extraction* step. ‘C’ of ‘0+1’ denotes only appending one linear layer after the module, and ‘D’ of ‘1+1’ denotes inserting

Similarity Metric	Acc (%)	Traditional Classifier	Acc (%)
Chebyshev	65.9	Dec. Tree [54]	57.8
Euclidean	79.8	GBoost [14]	63.9
Manhattan	80.9	SVM [8]	79.9
Cosine	81.8	Memory Bank	81.8

Table 9. **Ablation Study of Point-Memory Bank.** We ablate the metric for similarity calculation and compare our memory bank with traditional machine learning algorithms. ‘GBoost’ and ‘Dec. Tree’ denote GradientBoosting and Decision Tree algorithms.

Method	mIoU	Method	AP ₂₅	AR ₂₅
DGCNN	85.16	VoteNet	57.84	80.92
+NN	+0.16	+NN	+1.28	+5.31
CurveNet	86.58	3DETR-m	64.60	77.22
+NN	+0.07	+NN	+1.02	+11.05

Table 10. **Plug-and-play for Part Segmentation and 3D Object Detection.** We report the performance gain (%) on ShapeNetPart [76] and ScanNetV2 [9], respectively for the two tasks.

into both positions. The preceding layers pre-transform the point features to better reveal the local geometry, and the latter layers further parse the weighted features for high-level understanding. We observe that Point-PN of ‘1+2’ performs the best, and more learnable layers, ‘2+2’, would harm the performance.

4.4. Plug-and-play

Shape Classification. We evaluate the enhancement capacity of Point-NN on two classification datasets in Table 8. By simple inference-time ensemble, Point-NN effectively boosts existing methods with different margins. On the more challenging ScanObjectNN [62], both PointNet [45] and PointMLP [39] are improved by +2.0% accuracy. This well indicates the effectiveness of complementary knowledge provided by Point-NN.

Segmentation and Detection. In Table 10, we present the plug-and-play performance of Point-NN on part segmentation and 3D object detection tasks. As the segmentation scores have long been saturated on the ShapeNetPart [76] benchmark, the boost of +0.1% mIoU for state-of-the-art CurveNet [70] is still noteworthy. For detection, Point-NN significantly enhances 3DETR-m [42] by +1.02% AP₂₅ and +11.05% AR₂₅. By fusing complementary knowledge to the trained classifier, the 3D detector can better judge whether the candidate boxes include objects and correctly remove the false positive ones. The experiments fully indicate the general enhancement ability of Point-NN.

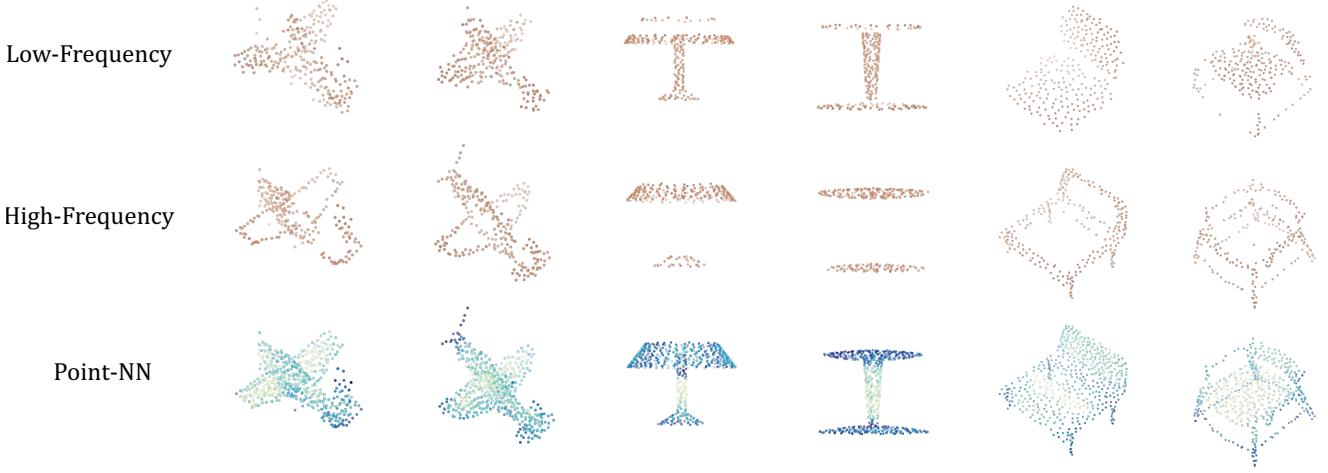


Figure 7. Why Do Trigonometric Functions Work? For an input point cloud, we visualize its low-frequency and high-frequency geometries referring to [72], and compare with the feature responses of Point-NN, where darker colors indicate higher responses. As shown, Point-NN can focus on the high-frequency 3D structures with sharp variations of the point cloud.

5. Discussion

5.1. Why Do Trigonometric Functions Work?

We leverage the trigonometric function to conduct non-parametric raw-point embedding and geometry extraction. It can reveal the 3D spatial patterns benefited from the following three properties.

Capturing High-frequency 3D Structures. As discussed in Tancik et al. [58], transforming low-dimensional input by sinusoidal mapping helps MLPs to learn the high-frequency content during training. Similarly to our non-parametric encoding, Point-NN utilizes trigonometric functions to capture the high-frequency spatial structures of 3D point clouds, and then recognize them from these distinctive characteristics by the point-memory bank. In Figure 7, we visualize the low-frequency (Top) and high-frequency (Middle) geometry of the input point cloud, and compare them with the feature responses of Point-NN (Bottom). The high-frequency geometries denotes the spatial regions of edges, corners, and other fine-grained details, where the local 3D coordinates vary dramatically, while the low-frequency structure normally includes some flat and smooth object surfaces with gentle variations. As shown, aided by trigonometric functions, our Point-NN can concentrate well on these high-frequency 3D patterns.

Encoding Absolute and Relative Positions. Benefited from the nature of sinusoid, the trigonometric functions can not only represent the absolute position in the embedding space, but also implicitly encode the relative positional information between two 3D points. For two points, $p_i = (x_i, y_i, z_i)$ and $p_j = (x_j, y_j, z_j)$, we first obtain their

C -dimensional embeddings referring to Equation (5~7) in the main paper, formulated as

$$\text{PosE}(p_i) = \text{Concat}(f_i^x, f_i^y, f_i^z), \quad (15)$$

$$\text{PosE}(p_j) = \text{Concat}(f_j^x, f_j^y, f_j^z), \quad (16)$$

where $\text{PosE}(\cdot)$ denotes the positional encoding by trigonometric functions, and $f_{i/j}^x, f_{i/j}^y, f_{i/j}^z \in \mathbb{R}^{1 \times \frac{C}{3}}$ denote the embeddings of three axes. Then, their spatial relative relation can be revealed by the dot production between the two embeddings, formulated as

$$f_i^x f_j^{xT} + f_i^y f_j^{yT} + f_i^z f_j^{zT} = \text{PosE}(p_i) \text{PosE}(p_j)^T.$$

Taking the x axis as an example,

$$\sum_{m=0}^{\frac{C}{6}-1} \cos(\alpha(x_i - x_j)/\beta^{\frac{6m}{C}}) = f_i^x f_j^{xT}, \quad (17)$$

which indicates the relative x-axis distance of two points, in a similar way to the other two axes. Therefore, the trigonometric function is capable of encoding both absolute and relative 3D positional information for point cloud analysis.

Local Geometry Extraction. With the analysis above, we can reinterpret the Equation (9) of reweighting in the main paper from a mathematical perspective, which can be formulated as

$$\begin{aligned} f_{cj}^w &= (f_{cj} + \text{PosE}(\Delta p_j)) \odot \text{PosE}(\Delta p_j) \\ &= f_{cj} \odot \text{PosE}(\Delta p_j) + \text{PosE}(\Delta p_j) \odot \text{PosE}(\Delta p_j), \end{aligned} \quad (18)$$

where \odot denotes element-wise multiplication, and $\Delta p_j, f_{cj}$ denote the relative coordinates and feature of neighbor point

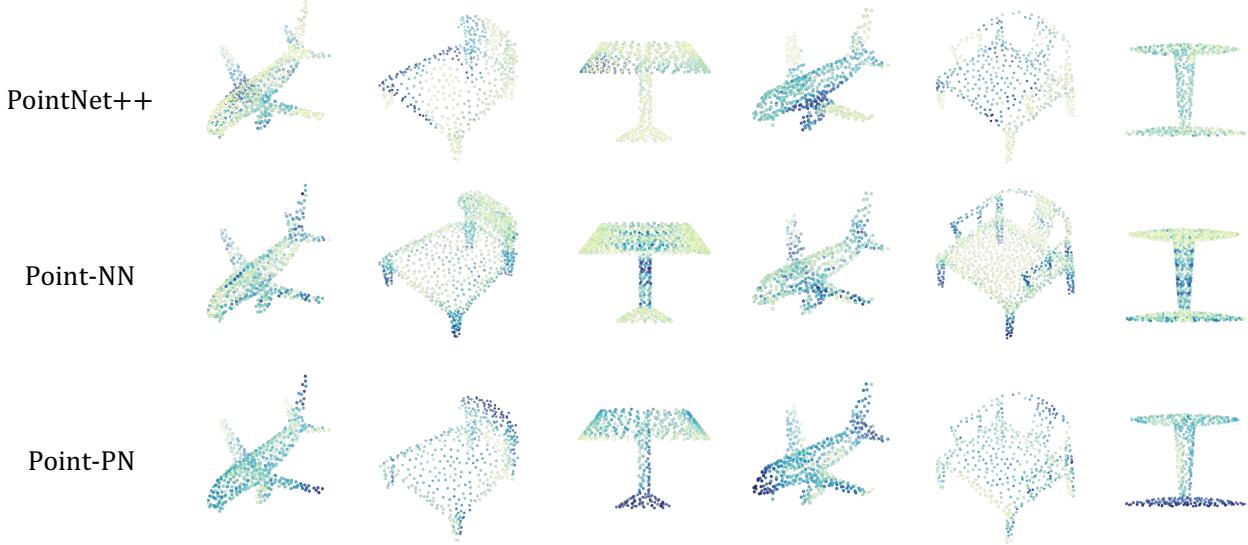


Figure 8. Can Point-NN Improve Point-PN by Plug-and-play? We visualize the feature responses for Point-NN, the trained PointNet++ [46] and Point-PN, where darker colors indicate higher responses. As shown, Point-PN captures similar 3D patterns to Point-NN, which harms their complementarity.

Benchmark	ScanObjectNN	ModelNet40	ShapeNetPart
Point-PN	87.1	93.8	86.6
+NN	+0.1	+0.2	+0.0

Table 11. Can Point-NN Improve Point-PN by Plug-and-play?
We report the accuracy (%) on the PB-T50-RS split of ScanObjectNN [62], ModelNet40 [69], and ShapeNetPart [76].

j to the center point c . Now, there are two terms in Equation 18, both of which reveal the relative geometry between point i and j . The former implicitly captures the relative positional information from point j to i , since the summation of its elements corresponds to the dot production between embeddings. The latter also indicates the embedding-space distance between the two points, which can be reformulated as the vector length of $\text{Pose}(\Delta p_j)$. In this way, Point-NN can effectively extract local 3D patterns and aggregate them by pooling operations.

5.2. Can Point-NN Improve Point-PN?

Point-NN can provide complementary geometric knowledge and serve as a plug-and-play module to boost existing learnable 3D models. Although Point-PN is also a learnable 3D network, the enhanced performance brought by Point-NN is marginal as reported in Table 11. By visualizing feature responses in Figure 8, we observe that the complementarity between Point-NN and Point-PN is much weaker than that between Point-NN and PointNet++ [46]. This is because the non-parametric framework of Point-PN is mostly

Method	Pre-train in 2D	Pre-train in 3D	3D Data	Acc. (%)
PointCLIP [79]	✓	-	-	20.2
CALIP [21]	✓	-	-	21.5
CLIP2Point [27]	✓	✓	✓	49.4
ULIP [74]	✓	✓	✓	60.4
PointCLIP V2 [91]	✓	-	-	64.2
Point-NN	-	-	✓	81.8

Table 12. Comparison of Training-free Methods in 3D. We report their performance without training on ModelNet40 [69].

inherited from Point-NN, also capturing high-frequency 3D geometries via trigonometric functions. Therefore, the learnable Point-PN extracts similar 3D patterns to Point-NN, which harms its plug-and-play capacity.

5.3. Training-free Methods in 3D

Our Point-NN conducts no training, but requires 3D training data to construct the point-memory bank. Inspired by the transfer learning in 2D and language [5, 15, 17, 34, 77, 80, 85, 90], some recent works [21, 27, 74, 79, 91] adapt the pre-trained models from other modalities, e.g., CLIP [50], to 3D domains in a zero-shot manner. Via the diverse pre-trained knowledge, they are also training-free and do not need any 3D training data. As compared in Table 12, different from other methods based on 2D or 3D pre-training, our method is a pure non-parametric network without any learnable parameters, which leverages feature similarities for recognition, other than pre-trained knowledge.

k	1	10	100	500	1000	5000	All
Top- k PoM	80.4	81.1	81.3	81.4	81.4	81.7	81.8
k -NN	80.7	79.5	67.0	45.7	36.4	8.5	-

Table 13. **Point-Memory Bank vs. k -NN.** ‘Top- k PoM’ denotes the point-memory bank with top- k similarities, and ‘All’ denotes 9,840 training samples. We utilize our non-parametric encoder to extract features and report the accuracy (%) on ModelNet40 [69].

5.4. Point-Memory Bank vs. k -NN?

Based on the already extracted point cloud features, our point-memory bank and k -NN algorithm both leverage the inter-sample feature similarity for classification, and require no training or learnable parameters, but are different from the following two aspects.

Soft Integration vs. Hard Assignment. As illustrated in Section (2.3) of the main paper, our point-memory bank regards the similarities S_{cos} between the test point cloud feature and the feature memory, F_{mem} , as weights, which are adopted for weighted summation of the one-hot label memory, T_{mem} . This can be viewed as a soft label integration. Instead, k -NN utilizes S_{cos} to search the k nearest neighbors from the training set, and directly outputs the category label with the maximum number of samples within the k neighbors. Hence, k -NN conducts a hard label assignment, which is less adaptive than the soft integration. Additionally, our point-memory bank can be accomplished simply by two matrix multiplications and requires no sorting, which is more efficient for hardware implementation.

All Samples vs. k Neighbors. Our point-memory bank softly integrates the entire label memory with different weights, i.e., similarity scores. This can take the semantics of all training samples into account for classification. In contrast, k -NN only involves the nearest k neighbors to the test sample, which discards the sufficient category knowledge from other training samples.

Performance Comparison. In Table 13, based on the point cloud features extracted by our non-parametric encoder, we implement the top- k version of point-memory bank for comparison with k -NN, which only aggregates the label memory of the training samples with top- k similarities. As the neighbor number k increases, k -NN’s performance is severely harmed due to its hard label assignment, while our point-memory bank attains the highest accuracy by utilizing all 9,840 samples for classification, indicating their different characters.

6. Conclusion

We revisit the non-learnable components in existing 3D models and propose Point-NN, a pure non-parametric network for 3D point cloud analysis. Free from any parameters or training, Point-NN achieves favorable accuracy on various 3D tasks. Starting from Point-NN, we propose its two promising applications: architectural frameworks for Point-PN and plug-and-play modules for performance improvement. Extensive experiments have demonstrated its effectiveness and significance. For future works, we will focus on exploring more advanced non-parametric models with wider application scenarios for 3D point cloud analysis.

A. Related Work

3D Point Cloud Analysis. As the main data form in 3D, point clouds have stimulated a range of challenging tasks, including shape classification [37, 39, 45–47], scene segmentation [10, 31, 88], 3D object detection [11, 23, 26, 42, 56, 81], and 3D vision-language learning [21, 68, 79, 91]. Existing solutions as backbone networks can be categorized into projection-based and point-based methods. To handle the irregularity and sparsity of point clouds, projection-based methods convert them into grid-like data, such as tangent planes [59], multi-view depth maps [18, 22, 57, 79, 91], and 3D voxels [38, 40, 52, 84]. By doing this, the efficient 2D networks [24] and 3D convolutions [38] can be adopted for robust point cloud understanding. However, the projection process inevitably causes geometric information loss and quantization error. Point-based methods directly extract 3D patterns upon the unstructured input points to alleviate this loss of information. The seminal PointNet [45] utilizes shared MLP layers to independently extract point features and aggregate the global representation via a max pooling. PointNet++ [46] further constructs a multi-stage hierarchy to encode local spatial geometries progressively. Since then, the follow-up methods introduce advanced yet complicated local operators [39, 71] and global transformers [19, 20, 78, 83] for spatial geometry learning. In this paper, we follow the paradigm of more popular point-based methods, and propose a pure non-parametric network, Point-NN, with its two promising applications. For the first time, we verify the effectiveness of non-parametric components for 3D point cloud analysis.

Local Geometry Operators. Referring to the inductive bias of locality [24, 30], most existing 3D models adopt delicate 3D operators to iteratively aggregate neighborhood features. Following PointNet++ [46], a series of methods utilize shared MLP layers with learnable relation modules for local pattern encoding, e.g., fully-linked webs [86], structural relation network [13], and geometric affine module [39]. Some methods define irregular spatial kernels and introduce point-wise convolutions by relation mapping [36], Monte Carlo estimation [25, 67], and dynamic kernel assembling [71]. Inspired by graph networks, DGCNN [65] and others [32, 60] regard points as vertices and interact local geometry through edges. Transformers [31, 87] are also introduced in 3D for attention-based feature communication. CurveNet [70] proposes generating hypothetical curves for point grouping and feature aggregation. Unlike all previous methods with learnable operators, Point-NN adopts non-parametric trigonometric functions to reveal the spatial geometry within local regions, and Point-PN further appends simple linear layers on top with high performance-parameter trade-off.

Positional Encodings. Transformers [63] represent input signals as an orderless sequence and implicitly utilize positional encodings (PE) to inject positional information. Typically, trigonometric functions are adopted as the non-learnable PE [16] to encode both absolute and relative positions, each dimension corresponding to a sinusoid. For vision, PE can also be learnable during training [12] or online predicted by neural networks [6, 87]. Another work [51] indicates that deep networks can learn better high-frequency variation given a higher dimensional input. Tancik *et al.* [58] interpret it as Fourier transform to learn high-frequency functions in low dimensional problems. NeRF [41] utilizes trigonometric PE to enhance the MLPs for better neural scene representations, but in a different formulation from the Transformer’s. In contrast, we extend the non-learnable trigonometric PE of Transformer for specialized raw-point embedding and local geometry extraction, other than serving as an accessory in previous learnable networks. By doing this, the non-parametric encoder of Point-NN can effectively capture low-level spatial patterns complementary to the already trained 3D models.

B. Implementation Details

Point-NN. The non-parametric encoder of Point-NN contains 4 stages. Each stage reduces the point number by half via FPS, and doubles the feature dimension during feature expansion. For shape classification, the initial feature dimension C_I is set to 72, and the final dimension C_G of global representation is 1,152. The neighbor number k of k -NN is 90 for all stages. We set the two hyperparameters α, β in PosE(\cdot) as 1000 and 100, respectively, referring to Equation (6) and (7) in the main paper. For part segmentation, we extend the non-parametric encoder into 5 stages for fully aggregate multi-scale 3D features. We set the initial feature dimension C_I as 144, and the neighbor number k as 128. We appended a non-parametric decoder with skip connections in each stage, which concatenate the propagated point features in the decoder with the corresponding ones from the encoder. As shown in Figure 9, we construct the segmentation point-memory bank by storing the part-wise features and labels from the training set, which largely saves the GPU memory. During inference, each point-wise feature of the test point cloud conducts similarity matching with the part-wise feature memory for segmentation.

Point-PN. For the parametric variant, we decrease C_I to 36 and the neighbor number k to 40 for lightweight parameters and efficient inference. We adopt the bottleneck architecture with a ratio 0.5 for the two cascaded linear layers after the *Geometry Extraction* step. The initial parametric raw-point embedding consists of only one linear layer, and the final classifier contains three linear layers as existing

Point-Memory Bank for Part Segmentation

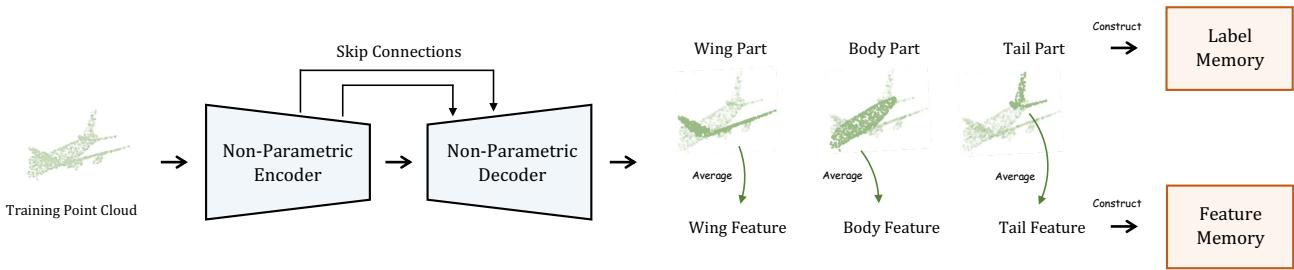


Figure 9. **Point-Memory Bank for Part Segmentation.** We first utilize the non-parametric encoder and decoder to extract the point-wise features of the point cloud in the training set. Then, we average the point features with the same part label to obtain the part-wise features, and construct them as the feature memory.

Grouping Method	Acc.	Feature Expand	Acc.	Pooling	Acc.
Ball Query	78.5	w/o	70.0	Max	80.4
<i>k</i> -NN	81.8	w	81.8	Ave	77.1
				Both	81.8

Table 14. **Ablation Study of Non-Parametric Encoder.** We ablate different designs at every stage of the encoder: the grouping method for local neighbors, feature expansion by concatenation, and pooling operation for feature aggregation. We report the classification accuracy (%) on ModelNet40 [69].

methods [39, 47]. For shape classification, we train Point-PN for 300 epochs with a batch size 32 on a single RTX 3090 GPU. On ModelNet40 [69], we adopt SGD optimizer with a weight decay 0.0002, and cosine scheduler with an initial learning rate 0.1. On ScanObjectNN [62], we adopt AdamW optimizer [29] with a weight decay 0.05, and cosine scheduler with an initial learning rate 0.002. We follow the data augmentation in [39] and [47] respectively for ModelNet40 and ScenObjectNN datasets. For part segmentation, we simply utilize the learnable decoder and training settings of CurveNet [70] for a fair comparison.

C. Additional Ablation Study

Non-Parametric Encoder. In Table 14, we further investigate other designs at every stage of Point-NN’s non-parametric encoder. As shown, *k*-NN performs better than ball query [46] for grouping the neighbors of each center point since the ball query would fail to aggregate valid geometry in some sparse regions with only a few neighboring points. Expanding the feature dimension by concatenating the center and neighboring points can improve the performance by +5.3%. This is because each point obtains larger receptive fields as the network stage goes deeper and requires higher-dimensional vectors to encode more spatial semantics. For the pooling operation after geometry extraction, we observe applying both max and average pooling

Magnitude α	1	10	50	100	200	500
Acc. (%)	68.3	77.9	81.1	81.8	81.4	81.0

Table 15. **Magnitude α in Trigonometric Functions.** We report the classification accuracy of Point-NN on ModelNet40 [69].

Wavelength β	10	100	500	1000	2000	3000
Acc. (%)	51.3	80.4	81.8	74.5	73.1	72.9

Table 16. **Wavelength β in Trigonometric Functions.** We report the classification accuracy of Point-NN on ModelNet40 [69].

achieves the highest accuracy, which can summarize the local patterns from two different aspects.

Hyperparameters in Trigonometric Functions. In Table 15 and 16, we show the influence of two hyperparameters in trigonometric functions of Point-NN. We fix one of them to be the best-performing value (α as 100, β as 500), and vary the other one for ablation. The combination of the magnitude α and wavelength β control the frequency of the channel-wise sinusoid, and thus determine the raw point encoding for different classification accuracy.

Point-Memory Bank with Different Sizes. As default, we construct the feature memory by the entire training-set point clouds. In Table 17, we report how Point-NN performs when partial training samples are utilized for the point-memory bank. As shown, Point-NN can attain 60.1% classification accuracy with only 10% of the training data, and further achieves 70.1% with 40% data, which is comparable to the performance of 100% ratio but consumes less GPU memory. This indicates Point-NN is not sensitive to the memory bank size and can perform favorably with partial training-set data.

Ratio (%)	1	5	10	20	40	80	100
Acc. (%)	39.5	64.2	70.3	75.0	77.9	80.8	81.8
Mem. (G)	3.84	3.87	3.93	4.05	4.26	4.82	5.21

Table 17. **Point-Memory Bank with Different Sizes.** We randomly sample different ratios of ModelNet40 [69] to construct the point-memory bank and report the classification accuracy with GPU memory consumption.

References

- [1] Aitor Aldoma, Zoltan-Csaba Marton, Federico Tombari, Walter Wohlkinger, Christian Potthast, Bernhard Zeisl, Radu Bogdan Rusu, Suat Gedikli, and Markus Vincze. Tutorial: Point cloud library: Three-dimensional object recognition and 6 dof pose estimation. *IEEE Robotics & Automation Magazine*, 19(3):80–91, 2012. [1](#)
- [2] Yizhak Ben-Shabat, Michael Lindenbaum, and Anath Fischer. 3dmfv: Three-dimensional point cloud classification in real-time using convolutional neural networks. *IEEE Robotics and Automation Letters*, 3(4):3145–3152, 2018. [2](#), [6](#), [7](#)
- [3] Jingdao Chen, Zsolt Kira, and Yong K Cho. Deep learning approach to point cloud scene understanding for automated scan to 3d reconstruction. *Journal of Computing in Civil Engineering*, 33(4):04019027, 2019. [1](#)
- [4] Xiaozhi Chen, Huimin Ma, Ji Wan, Bo Li, and Tian Xia. Multi-view 3d object detection network for autonomous driving. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1907–1915, 2017. [1](#)
- [5] Yuxiao Chen, Jianbo Yuan, Yu Tian, Shijie Geng, Xinyu Li, Ding Zhou, Dimitris N. Metaxas, and Hongxia Yang. Revisiting multimodal representation in contrastive learning: from patch and token embeddings to finite discrete tokens. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. [10](#)
- [6] Xiangxiang Chu, Zhi Tian, Bo Zhang, Xinlong Wang, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Conditionalpositional encodings for vision transformers. *arXiv preprint arXiv:2102.10882*, 2021. [12](#)
- [7] Nikolaus Correll, Kostas E Bekris, Dmitry Berenson, Oliver Brock, Albert Causo, Kris Hauser, Kei Okada, Alberto Rodriguez, Joseph M Romano, and Peter R Wurman. Analysis and observations from the first amazon picking challenge. *IEEE Transactions on Automation Science and Engineering*, 15(1):172–188, 2016. [1](#)
- [8] Nello Cristianini, John Shawe-Taylor, et al. *An introduction to support vector machines and other kernel-based learning methods*. Cambridge university press, 2000. [8](#)
- [9] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5828–5839, 2017. [2](#), [7](#), [8](#)
- [10] Angela Dai and Matthias Nießner. 3dmv: Joint 3d-multi-view prediction for 3d semantic scene segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 452–468, 2018. [12](#)
- [11] Zhipeng Ding, Xu Han, and Marc Niethammer. Votenet: A deep learning label fusion method for multi-atlas segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 202–210. Springer, 2019. [5](#), [7](#), [12](#)
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [12](#)
- [13] Yueqi Duan, Yu Zheng, Jiwen Lu, Jie Zhou, and Qi Tian. Structural relational reasoning of point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 949–958, 2019. [12](#)
- [14] Jerome H Friedman. Greedy function approximation: a gradient boosting machine. *Annals of statistics*, pages 1189–1232, 2001. [8](#)
- [15] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and Yu Qiao. Clip-adapter: Better vision-language models with feature adapters. *arXiv preprint arXiv:2110.04544*, 2021. [10](#)
- [16] Jonas Gehring, Michael Auli, David Grangier, Denis Yarats, and Yann N Dauphin. Convolutional sequence to sequence learning. In *International conference on machine learning*, pages 1243–1252. PMLR, 2017. [12](#)
- [17] Shijie Geng, Jianbo Yuan, Yu Tian, Yuxiao Chen, and Yongfeng Zhang. HiCLIP: Contrastive language-image pre-training with hierarchy-aware attention. In *The Eleventh International Conference on Learning Representations*, 2023. [10](#)
- [18] Ankit Goyal, Hei Law, Bowei Liu, Alejandro Newell, and Jia Deng. Revisiting point cloud shape classification with a simple and effective baseline. *arXiv preprint arXiv:2106.05304*, 2021. [6](#), [12](#)
- [19] Meng-Hao Guo, Jun-Xiong Cai, Zheng-Ning Liu, Tai-Jiang Mu, Ralph R Martin, and Shi-Min Hu. Pct: Point cloud transformer. *Computational Visual Media*, 7(2):187–199, 2021. [1](#), [6](#), [12](#)
- [20] Ziyu Guo, Xianzhi Li, and Pheng Ann Heng. Joint-mae: 2d-3d joint masked autoencoders for 3d point cloud pre-training. *arXiv preprint arXiv:2302.14007*, 2023. [12](#)
- [21] Ziyu Guo, Renrui Zhang, Longtian Qiu, Xianzheng Ma, Xupeng Miao, Xuming He, and Bin Cui. Calip: Zero-shot enhancement of clip with parameter-free attention. *arXiv preprint arXiv:2209.14169*, 2022. [10](#), [12](#)
- [22] Abdullah Hamdi, Silvio Giancola, and Bernard Ghanem. Mvtn: Multi-view transformation network for 3d shape recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1–11, 2021. [12](#)
- [23] Chenzhang He, Hui Zeng, Jianqiang Huang, Xian-Sheng Hua, and Lei Zhang. Structure aware single-stage 3d object detection from point cloud. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11873–11882, 2020. [12](#)

- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 12
- [25] Pedro Hermosilla, Tobias Ritschel, Pere-Pau Vázquez, Àlvar Vinacua, and Timo Ropinski. Monte carlo convolution for learning on non-uniformly sampled point clouds. *ACM Transactions on Graphics (TOG)*, 37(6):1–12, 2018. 12
- [26] Peixiang Huang, Li Liu, Renrui Zhang, Song Zhang, Xinli Xu, Baichao Wang, and Guoyi Liu. Tig-bev: Multi-view bev 3d object detection via target inner-geometry learning. *arXiv preprint arXiv:2212.13979*, 2022. 12
- [27] Tianyu Huang, Bowen Dong, Yunhan Yang, Xiaoshui Huang, Rynson WH Lau, Wanli Ouyang, and Wangmeng Zuo. Clip2point: Transfer clip to point cloud classification with image-depth pre-training. *arXiv preprint arXiv:2210.01055*, 2022. 10
- [28] Kiyosumi Kidono, Takeo Miyasaka, Akihiro Watanabe, Takashi Naito, and Jun Miura. Pedestrian recognition using high-definition lidar. In *2011 IEEE Intelligent Vehicles Symposium (IV)*, pages 405–410. IEEE, 2011. 1
- [29] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 13
- [30] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017. 12
- [31] Xin Lai, Jianhui Liu, Li Jiang, Liwei Wang, Hengshuang Zhao, Shu Liu, Xiaojuan Qi, and Jiaya Jia. Stratified transformer for 3d point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8500–8509, 2022. 12
- [32] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4558–4567, 2018. 12
- [33] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhuan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. *Advances in neural information processing systems*, 31:820–830, 2018. 2, 6, 7
- [34] Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard de Melo, Xiaogang Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXV*, pages 388–404. Springer, 2022. 10
- [35] Yongcheng Liu, Bin Fan, Gaofeng Meng, Jiwen Lu, Shiming Xiang, and Chunhong Pan. Densepoint: Learning densely contextual representation for efficient point cloud processing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5239–5248, 2019. 1, 6
- [36] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8895–8904, 2019. 6, 12
- [37] Ze Liu, Han Hu, Yue Cao, Zheng Zhang, and Xin Tong. A closer look at local aggregation operators in point cloud analysis. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 326–342. Springer, 2020. 12
- [38] Zhijian Liu, Haotian Tang, Yujun Lin, and Song Han. Point-voxel cnn for efficient 3d deep learning. *arXiv preprint arXiv:1907.03739*, 2019. 1, 12
- [39] Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local geometry in point cloud: A simple residual mlp framework. *arXiv preprint arXiv:2202.07123*, 2022. 1, 2, 6, 8, 12, 13
- [40] Hsien-Yu Meng, Lin Gao, Yu-Kun Lai, and Dinesh Manocha. Vv-net: Voxel vae net with group convolutions for point cloud segmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 8500–8508, 2019. 12
- [41] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 8, 12
- [42] Ishan Misra, Rohit Girdhar, and Armand Joulin. An end-to-end transformer model for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 2906–2917, October 2021. 5, 7, 8, 12
- [43] Arsalan Mousavian, Clemens Eppner, and Dieter Fox. 6-dof graspernet: Variational grasp generation for object manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2901–2910, 2019. 1
- [44] Luis E Navarro-Serment, Christoph Mertz, and Martial Hebert. Pedestrian detection and tracking using three-dimensional lidar data. *The International Journal of Robotics Research*, 29(12):1516–1528, 2010. 1
- [45] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas 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*, pages 652–660, 2017. 2, 6, 7, 8, 12
- [46] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *arXiv preprint arXiv:1706.02413*, 2017. 1, 2, 5, 6, 7, 10, 12, 13
- [47] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Abed Al Kader Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies. *arXiv preprint arXiv:2206.04670*, 2022. 12, 13
- [48] Shi Qiu, Saeed Anwar, and Nick Barnes. Dense-resolution network for point cloud classification and segmentation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3813–3822, 2021. 6
- [49] Shi Qiu, Saeed Anwar, and Nick Barnes. Geometric back-projection network for point cloud classification. *IEEE Transactions on Multimedia*, 2021. 6
- [50] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry,

- Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 10
- [51] Nasim Rahaman, Aristide Baratin, Devansh Arpit, Felix Draxler, Min Lin, Fred Hamprecht, Yoshua Bengio, and Aaron Courville. On the spectral bias of neural networks. In *International Conference on Machine Learning*, pages 5301–5310. PMLR, 2019. 12
- [52] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3577–3586, 2017. 12
- [53] Radu Bogdan Rusu, Nico Blodow, Zoltan Csaba Marton, and Michael Beetz. Close-range scene segmentation and reconstruction of 3d point cloud maps for mobile manipulation in domestic environments. In *2009 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1–6. IEEE, 2009. 1
- [54] S Rasoul Safavian and David Landgrebe. A survey of decision tree classifier methodology. *IEEE transactions on systems, man, and cybernetics*, 21(3):660–674, 1991. 8
- [55] Charu Sharma and Manohar Kaul. Self-supervised few-shot learning on point clouds. *arXiv preprint arXiv:2009.14168*, 2020. 7
- [56] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. Pv-rcnn: Point-voxel feature set abstraction for 3d object detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020. 12
- [57] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015. 12
- [58] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems*, 33:7537–7547, 2020. 8, 9, 12
- [59] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3887–3896, 2018. 12
- [60] Gusi Te, Wei Hu, Amin Zheng, and Zongming Guo. Rgcn: Regularized graph cnn for point cloud segmentation. In *Proceedings of the 26th ACM international conference on Multimedia*, pages 746–754, 2018. 12
- [61] Hugues Thomas, Charles R Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6411–6420, 2019. 1
- [62] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1588–1597, 2019. 2, 6, 7, 8, 10, 13
- [63] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 3, 7, 12
- [64] Francesco Verdoja, Diego Thomas, and Akihiro Sugimoto. Fast 3d point cloud segmentation using supervoxels with geometry and color for 3d scene understanding. In *2017 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1285–1290. IEEE, 2017. 1
- [65] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *AcM Transactions On Graphics (tog)*, 38(5):1–12, 2019. 1, 6, 7, 12
- [66] Jiajun Wu, Chengkai Zhang, Tianfan Xue, William T Freeman, and Joshua B Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pages 82–90, 2016. 7
- [67] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019. 12
- [68] Yanmin Wu, Xinhua Cheng, Renrui Zhang, Zesen Cheng, and Jian Zhang. Eda: Explicit text-decoupling and dense alignment for 3d visual and language learning. *arXiv preprint arXiv:2209.14941*, 2022. 12
- [69] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Liguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1912–1920, 2015. 5, 6, 7, 8, 10, 11, 13, 14
- [70] Tiange Xiang, Chaoyi Zhang, Yang Song, Jianhui Yu, and Weidong Cai. Walk in the cloud: Learning curves for point clouds shape analysis. *arXiv preprint arXiv:2105.01288*, 2021. 1, 6, 8, 12, 13
- [71] Mutian Xu, Runyu Ding, Hengshuang Zhao, and Xiaojuan Qi. Paconv: Position adaptive convolution with dynamic kernel assembling on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3173–3182, 2021. 1, 12
- [72] Mutian Xu, Junhao Zhang, Zhipeng Zhou, Mingye Xu, Xiaojuan Qi, and Yu Qiao. Learning geometry-disentangled representation for complementary understanding of 3d object point cloud. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3056–3064, 2021. 9
- [73] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 87–102, 2018. 6
- [74] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles,

- and Silvio Savarese. Ulip: Learning unified representation of language, image and point cloud for 3d understanding. *arXiv preprint arXiv:2212.05171*, 2022. 10
- [75] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 206–215, 2018. 7
- [76] Li Yi, Vladimir G Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Transactions on Graphics (ToG)*, 35(6):1–12, 2016. 6, 7, 8, 10
- [77] Renrui Zhang, Rongyao Fang, Peng Gao, Wei Zhang, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *arXiv preprint arXiv:2111.03930*, 2021. 4, 10
- [78] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Pointm2ae: Multi-scale masked autoencoders for hierarchical point cloud pre-training. *arXiv preprint arXiv:2205.14401*, 2022. 12
- [79] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8552–8562, 2022. 10, 12
- [80] Renrui Zhang, Xiangfei Hu, Bohao Li, Siyuan Huang, Hanqiu Deng, Hongsheng Li, Yu Qiao, and Peng Gao. Prompt, generate, then cache: Cascade of foundation models makes strong few-shot learners. *arXiv preprint arXiv:2303.02151*, 2023. 10
- [81] Renrui Zhang, Han Qiu, Tai Wang, Xuanzhuo Xu, Ziyu Guo, Yu Qiao, Peng Gao, and Hongsheng Li. Monodetr: Depth-aware transformer for monocular 3d object detection. *arXiv preprint arXiv:2203.13310*, 2022. 12
- [82] Renrui Zhang, Liuhui Wang, Ziyu Guo, and Jianbo Shi. Nearest neighbors meet deep neural networks for point cloud analysis. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1246–1255, 2023. 6
- [83] Renrui Zhang, Liuhui Wang, Yu Qiao, Peng Gao, and Hongsheng Li. Learning 3d representations from 2d pre-trained models via image-to-point masked autoencoders. *arXiv preprint arXiv:2212.06785*, 2022. 12
- [84] Renrui Zhang, Ziyao Zeng, Ziyu Guo, Xinben Gao, Kexue Fu, and Jianbo Shi. Dspoint: Dual-scale point cloud recognition with high-frequency fusion. *arXiv preprint arXiv:2111.10332*, 2021. 12
- [85] Renrui Zhang, Ziyao Zeng, Ziyu Guo, and Yafeng Li. Can language understand depth? In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 6868–6874, 2022. 10
- [86] Hengshuang Zhao, Li Jiang, Chi-Wing Fu, and Jiaya Jia. Pointweb: Enhancing local neighborhood features for point cloud processing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5565–5573, 2019. 12
- [87] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021. 1, 12
- [88] Zongyue Zhao, Min Liu, and Karthik Ramani. Dar-net: Dynamic aggregation network for semantic scene segmentation. *arXiv preprint arXiv:1907.12022*, 2019. 12
- [89] Bo Zheng, Yibiao Zhao, Joey C Yu, Katsushi Ikeuchi, and Song-Chun Zhu. Beyond point clouds: Scene understanding by reasoning geometry and physics. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3127–3134, 2013. 1
- [90] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *International Journal of Computer Vision*, 130(9):2337–2348, 2022. 10
- [91] Xiangyang Zhu, Renrui Zhang, Bowei He, Ziyao Zeng, Shanghang Zhang, and Peng Gao. Pointclip v2: Adapting clip for powerful 3d open-world learning. *arXiv preprint arXiv:2211.11682*, 2022. 10, 12