

DensePoint: Learning Densely Contextual Representation for Efficient Point Cloud Processing

Yongcheng Liu^{†‡} Bin Fan^{*†} Gaofeng Meng[†] Jiwen Lu[§] Shiming Xiang^{†‡} Chunhong Pan[†]

[†]National Laboratory of Pattern Recognition, Institute of Automation, Chinese Academy of Sciences

[‡]School of Artificial Intelligence, University of Chinese Academy of Sciences

[§]Department of Automation, Tsinghua University

{yongcheng.liu, bfan, gfmeng, smxiang, chpan}@nlpr.ia.ac.cn lujiwen@tsinghua.edu.cn

Abstract

Point cloud processing is very challenging, as the diverse shapes formed by irregular points are often indistinguishable. A thorough grasp of the elusive shape requires sufficiently contextual semantic information, yet few works devote to this. Here we propose DensePoint, a general architecture to learn densely contextual representation for point cloud processing. Technically, it extends regular grid CNN to irregular point configuration by generalizing a convolution operator, which holds the permutation invariance of points, and achieves efficient inductive learning of local patterns. Architecturally, it finds inspiration from dense connection mode, to repeatedly aggregate multi-level and multi-scale semantics in a deep hierarchy. As a result, densely contextual information along with rich semantics, can be acquired by DensePoint in an organic manner, making it highly effective. Extensive experiments on challenging benchmarks across four tasks, as well as thorough model analysis, verify DensePoint achieves the state of the arts.

1. Introduction

Recently, the processing of point cloud, which comprises an irregular set of 3D points, has drawn a lot of attention, due to its wide range of applications such as robot manipulation [20] and autonomous driving [31]. However, modern applications usually demand for a high-level understanding of point cloud, *i.e.*, identifying the implicit 3D shape pattern. This is quite challenging, since the diverse shapes, abstractly formed by these irregular points, are often hardly distinguishable. For this issue, it is essential to capture sufficiently contextual semantic information for a thorough grasp of the elusive shape (see Fig. 1 for details).

Over the past few years, convolutional neural network (CNN) has demonstrated its powerful abstraction ability of semantic information in image recognition field [61]. Ac-

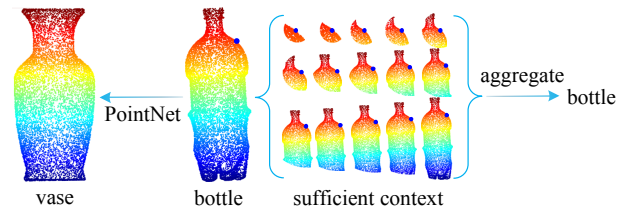


Figure 1. **Motivation:** sufficiently contextual semantic information is essential for a thorough grasp of the elusive shape formed by point cloud. The “bottle” is misidentified as the “vase” by PointNet [32], while with sufficient context aggregated, it can be accurately recognized. Here, we only illustrate the multi-level context around the blue point for visual clearness.

cordingly, much effort is focused on replicating its remarkable success on the analysis of image [22, 41], *i.e.*, regular grid data, to irregular point cloud processing [34, 19, 35, 49, 59, 27]. A straightforward strategy is to transform point cloud into regular voxels [53, 28, 4] or multi-view images [44, 3, 6], for easy application of CNN. These transformations, however, usually lead to much loss of rich 3D geometric information, as well as high complexity.

Another difficult yet attractive solution is to learn directly from irregular point cloud. PointNet [32], a pioneer in this direction, achieves the permutation invariance of points by learning over each point independently, then applying a symmetric function to accumulate features. Though impressive, it ignores local patterns that have been proven to be important for abstracting high-level visual semantics in image CNN [61]. To remedy this defect, KCNet [39] mines local patterns by creating a k-NN graph over each point in PointNet. Nevertheless, it inherits another defect of PointNet, *i.e.*, no pooling layer to explicitly raise the level of semantics. PointNet++ [34] hierarchically groups point cloud into local subsets and learns on them by PointNet. This design indeed works like CNN, but the basic operator, PointNet, demands high complexity for enough effectiveness.

Besides high-level semantics, contextual information, which reflects the potential semantic dependencies between a target pattern and its surroundings [30], is also critical for

*Corresponding author: Bin Fan

shape pattern recognition. A typical approach in this view is multi-scale learning. Accordingly, PointNet++ [34] directly applies multi-scale grouping in each layer, *i.e.*, capturing context at the same semantic level. This way, however, is suboptimal as it ignores the inherent difference in semantic levels at different scales, and often causes huge computational cost, especially for lots of scales. Multi-resolution grouping [34] can partly alleviate the latter issue, yet actually, it also abandons crucial context acquisition. ShapeContextNet [55] finds another strategy inspired by shape context [2]. It applies self-attention [48] in each layer of PointNet [32] to dynamically learn the relation weight among all points, and regards this weight as global shape context. Though fully automatic, it lacks an explicit semantic abstraction like CNN from local to global, and the weight matrix $N \times N$ in self-attention can cause huge complexity when the number of points N increases.

In short, there are mainly two key requirements to exploit CNN for effective learning on point cloud: 1) A convolution operator on point cloud, which can be permutation invariant to unordered points, and can achieve efficient inductive learning of local patterns, is required; 2) A deep hierarchy, which can acquire sufficiently contextual semantics for accurate shape recognition, is also required.

Accordingly, we propose DensePoint, a general architecture to learn densely contextual representation for point cloud processing, as illustrated in Fig. 2. Technically, DensePoint extends regular grid CNN to irregular point configuration by generalizing a convolution operator, which holds the permutation invariance of points, and respects the convolutional properties, *i.e.*, local connectivity and weight sharing. Owing to its efficient inductive learning of local patterns, a deep hierarchy can be easily built in DensePoint for semantic abstraction. Architecturally, DensePoint finds inspiration from dense connection mode [13], to repeatedly aggregate multi-level and multi-scale semantics in the deep hierarchy. As a result, densely contextual information along with rich semantics, can be acquired by DensePoint in an organic manner, making it highly effective.

The key contributions are highlighted as follows:

- A generalized convolution operator is formulated. It is permutation invariant to points, and respects the convolutional properties of local connectivity and weight sharing, thus extending regular grid CNN to irregular configuration for efficient point cloud processing.
- A general architecture equipped with the generalized convolution operator to learn densely contextual representation of point cloud, *i.e.*, DensePoint, is proposed. It can acquire sufficiently contextual semantic information for accurate recognition of the implicit shape.
- Comprehensive experiments on challenging benchmarks across four tasks, *i.e.*, shape classification, shape retrieval, part segmentation and normal estimation, as well as thorough model analysis, demonstrate that DensePoint achieves the state of the arts.

2. Related Work

In this section, we briefly review existing deep learning methods for 3D shape learning.

View-based and volumetric methods. View-based methods [44, 3, 6, 54, 7, 33, 14] represent a 3D shape as a collection of 2D views, over which classic CNN used in image analysis field can be easily applied. However, 2D projections could cause much loss of 3D shape information due to many self-occlusions. Volumetric methods convert a 3D shape into a regular 3D grid [53, 28, 4], over which 3D CNN [47] can be employed. The main limitation is the quantization loss of 3D shape information due to the low resolution enforced by 3D grid. Although this issue can be partly rescued by recent space partition methods like K-d trees [21] or octrees [50, 45, 36, 51], they still rely on a subdivision of a bounding volume. By contrast, our work devotes to learn directly from irregular 3D point cloud.

Deep learning on point cloud. Much effort has been focused on learning directly on point cloud. PointNet [32] pioneers this route by learning on each point independently and accumulating the final features. Yet it ignores local patterns, which limits its semantic learning ability. Accordingly, some works [34, 5, 39] partition point cloud into local subsets and learn on them based on PointNet. Some other works introduce graph convolutional network to learn over a local graph [49, 46, 25] or geometric elements [23]. However, these methods either lack an explicit semantic abstraction like CNN from local to global, or cause considerable complexity. By contrast, our work extends regular grid CNN to irregular point configuration, achieving efficient learning for point cloud processing.

In addition, there are some works mapping point cloud into a regular space to facilitate the application of classic CNN, *e.g.*, a sparse lattice structure [43] with bilateral convolution [18] or a continuous volumetric function [1] with 3D CNN. Nevertheless, in our case, we learn directly from irregular point cloud, which is much more challenging.

Contextual learning on point cloud. Contextual information is important for identifying the implicit shape pattern. PointNet++ [34] follows the traditional multi-scale learning by directly capturing context on the same layer, which often causes huge complexity. Hence an alternate called multi-resolution grouping [34] is devised for efficiency. It forces each layer to learn from its previous layer and the raw input (on the same local region) simultaneously. However, this can be less effective as it actually abandons crucial context acquisition. ShapeContextNet [55] finds another strategy inspired by shape context [2]. Instead of the traditional handcrafted design, it applies self-attention [48] to dynamically learn a weight for all point pairs. Though fully automatic, it lacks a local-to-global semantic learning like CNN. By contrast, we develop a deep hierarchy by an efficient generalized convolution operator, and organically aggregate multi-level contextual semantics in this hierarchy.

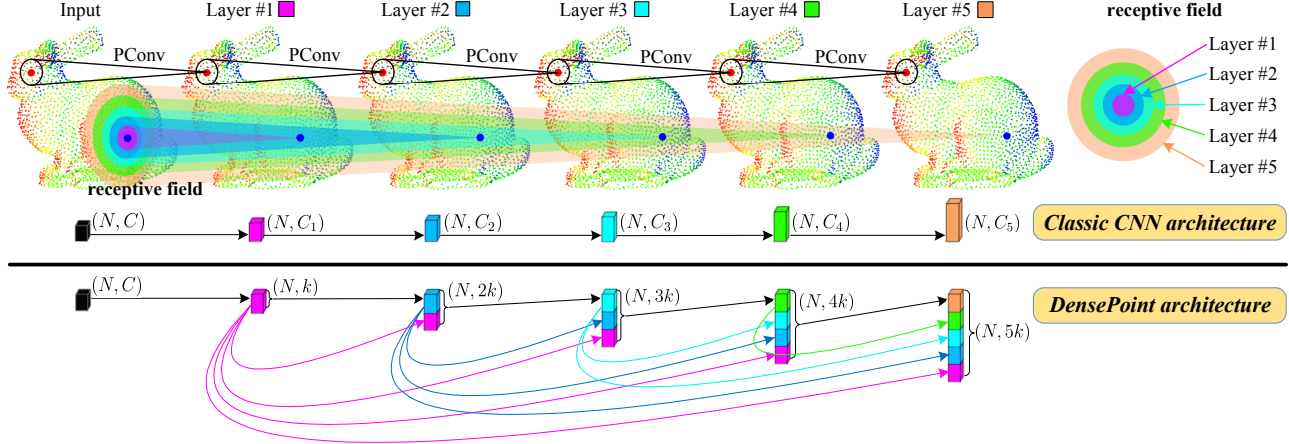


Figure 2. The illustration of DensePoint. It extends regular grid CNN to irregular point configuration by an efficient generalized convolution operator (PConv in Eq. (1)). Instead of classic CNN architecture with layer-by-layer connections, it finds inspiration from dense connection mode [13], to repeatedly aggregate multi-level along with multi-scale semantics in an organic manner. To avoid high complexity in deep layers, it forces the output of each layer to be equally narrow with a small constant k (e.g., 24). As a result, densely contextual representation can be learned efficiently for point cloud processing. Here, N is the number of points while C , C_* and k denote feature dimension.

3. Method

In this section, we first describe the generalized convolution operator and the pooling operator on point cloud. Then, we present DensePoint, and elaborate how it learns densely contextual representation for point cloud processing.

3.1. Convolution and Pooling on Point Cloud

PConv: convolution on point cloud. Classic convolution on the image operates on a local grid region (*i.e.*, local connectivity), and the convolution filter weights of this grid region are shared along the spatial dimension (*i.e.*, weight sharing). However, this operation is difficult to implement on point cloud due to the irregularity of points. To deal with this problem, we decompose the classic convolution into two core steps, *i.e.*, feature transformation and feature aggregation. Accordingly, a generalized convolution on point cloud can be formulated as

$$\mathbf{f}_{\mathcal{N}(x)} = \rho(\{\phi(\mathbf{f}_{x_n}), \forall x_n \in \mathcal{N}(x)\}), \quad (1)$$

where both x and x_n denote a 3D point in \mathbb{R}^3 , and \mathbf{f} is feature vector. $\mathcal{N}(x)$, the neighborhood formed by a local point cloud to convolve, is sampled from the whole point cloud by taking a sampled point x as the centroid, and the nearby points as its neighbors x_n . $\mathbf{f}_{\mathcal{N}(x)}$, the convolutional result as the inductive representation of $\mathcal{N}(x)$, is obtained by: (i) performing a feature transformation with function ϕ on each point in $\mathcal{N}(x)$; (ii) applying an aggregation function ρ to aggregate these transformed features. Finally, as shown in the upper part of Fig. 2 (PConv), similar to classic grid convolution, $\mathbf{f}_{\mathcal{N}(x)}$ is assigned to be the feature vector of the centroid point x in the next layer. Noticeably, some previous works such as [34] also use this general formulation.

In Eq. (1), $\mathbf{f}_{\mathcal{N}(x)}$ can be permutation invariant only when the inner function ϕ is shared over each point in $\mathcal{N}(x)$, and

the outer function ρ is symmetric (*e.g.*, sum). Accordingly, for high efficiency, we employ a shared single-layer perceptron (SLP, for short) following a nonlinear activator, as ϕ to implement feature transformation. Meanwhile, as done in classic convolution, ϕ is also shared over each local neighborhood, for achieving the weight sharing mechanism. As a result, with a symmetric ρ , the generalized PConv can achieve efficient inductive learning of local patterns, whilst be independent of the irregularity of points. Further, using PConv as the basic operator, a classic CNN architecture (no downsampling), as shown in the upper part of Fig. 2, can be easily built with layer-by-layer connections.

PPool: pooling on point cloud. In classic CNN, pooling is usually performed to explicitly raise the semantic level of the representation and improve computational efficiency. Here, using PConv, this operation can be achieved on point cloud in a learnable way. Specifically, N_o points are first uniformly sampled from the input N_i points, where $N_o < N_i$ (*e.g.*, $N_o = N_i/2$). Then, PConv can be applied to convolve all the local neighborhoods centered on those N_o points, to generate a new downsampling layer.

3.2. Learning Densely Contextual Representation

Classic CNN architecture. In a classic CNN architecture with layer-by-layer connections (the upper part of Fig. 2), hierarchical representations can be learned with the low-level ones in early layers and the high-level ones in deep layers [61]. However, a significant drawback is that each layer can only learn from single-level representation. As a consequence, all layers can capture only single-scale shape information from the input point cloud. Formally, assume a point cloud \mathbf{P}^0 that is passed through this type of network. The network comprises L layers, in which the ℓ^{th} layer performs a non-linear transformation $\mathcal{H}^\ell(\cdot)$. Then, the output

of the ℓ^{th} layer can be learned from its previous layer as

$$\mathbf{P}^\ell = \mathcal{H}^\ell(\mathbf{P}^{\ell-1}), \quad (2)$$

where each point in $\mathbf{P}^{\ell-1}$ is of single-scale receptive field on the input point cloud \mathbf{P}^0 , resulting that the learned \mathbf{P}^ℓ captures only single-scale shape information. Finally, this will lead to a weakly contextual representation, which is not effective enough for identifying the diverse implicit shapes.

DensePoint architecture. To overcome the above issue, we present a general architecture, *i.e.*, DensePoint shown in the lower part of Fig. 2, inspired by dense connection mode [13]. Specifically, for each layer in DensePoint (no downsampling), the outputs of all preceding layers are used as its input, and its own output is used as the input to all subsequent layers. That is, \mathbf{P}^ℓ in Eq. (2) becomes

$$\mathbf{P}^\ell = \mathcal{H}^\ell([\mathbf{P}^0, \mathbf{P}^1, \dots, \mathbf{P}^{\ell-1}]), \quad (3)$$

where $[\cdot]$ denotes the concatenation of the outputs of all preceding layers. Here, \mathbf{P}^ℓ is forced to learn from multi-level representations, which facilitates to aggregate multi-level shape semantics along with multi-scale shape information. In this way, each layer in DensePoint can capture a certain level (scale) of context, and the level can be gradually increased as the network deepens. Moreover, the acquired dense context in deep layers can also improve the abstraction of high-level semantics in turn, making the whole learning process organic. Eventually, very rich local-to-global shape information in the input \mathbf{P}^0 can be *progressively* aggregated together, resulting in a densely contextual representation for point cloud processing.

Note that DensePoint is quite different from the traditional multi-scale strategy [34]. The former progressively aggregates multi-level (multi-scale) semantics that is organically learned by each layer, while the latter artlessly gathers multi-scale information at the same level. It is also dissimilar to a simple concatenation of all layers as the final output, which results in each layer being less contextual.

Narrow architecture. When the network deepens, DensePoint will suffer from high complexity, since the convolutional overhead of deep layers will be huge with all preceding layers as the input. Thus, we narrow the output channels of each layer in DensePoint with a small constant k (*e.g.*, 24), instead of the large ones (*e.g.*, 512) in classic CNN.

ePConv: enhanced PConv. Though lightweight, such narrow DensePoint will lack the expressive power, since with much narrow output k , the shared SLP in PConv, *i.e.*, ϕ in Eq. (1), could be insufficient in terms of learning ability. To overcome this issue, we introduce the filter grouping [22] to enhance PConv, which divides all the filters in a layer into several groups, and each group performs individual operation. Formally, the enhanced PConv (ePConv, for short) converts Eq. (1) to

$$\mathbf{f}_{\mathcal{N}(x)} = \psi(\rho(\{\tilde{\phi}(\mathbf{f}_{x_n}), \forall x_n \in \mathcal{N}(x)\})), \quad (4)$$

Algorithm 1: DensePoint forward pass algorithm

Input: point cloud \mathbf{P} ; input features $\{\mathbf{f}_x, \forall x \in \mathbf{P}\}$; depth L ; weight $\mathbf{W}_{\tilde{\phi}}^\ell$, \mathbf{W}_ψ^ℓ and bias $\mathbf{b}_{\tilde{\phi}}^\ell$, \mathbf{b}_ψ^ℓ for SLP $\tilde{\phi}$ and SLP ψ in Eq. (4), $\forall \ell \in \{1, \dots, L\}$; non-linearity σ ; aggregation function ρ ; neighborhood method \mathcal{N}

Output: densely contextual representations $\{\mathbf{r}_x, \forall x \in \mathbf{P}\}$

- 1 $\mathbf{f}_x^0 \leftarrow \mathbf{f}_x, \forall x \in \mathbf{P}$;
- 2 **for** $\ell = 1 \dots L$ **do**
- 3 **for** $x \in \mathbf{P}$ **do**
- 4 $\mathbf{f}_{\mathcal{N}(x)}^\ell \leftarrow \rho(\{\sigma(\tilde{\mathbf{W}}_{\tilde{\phi}}^\ell \cdot \tilde{\mathbf{f}}_{x_n}^{\ell-1} + \tilde{\mathbf{b}}_{\tilde{\phi}}^\ell), \forall x_n \in \mathcal{N}(x)\})$;
- 5 $\mathbf{f}_x^\ell \leftarrow \sigma(\mathbf{W}_\psi^\ell \cdot \mathbf{f}_{\mathcal{N}(x)}^\ell + \mathbf{b}_\psi^\ell)$;
- 6 **end**
- 7 $\mathbf{f}_x^\ell \leftarrow [\mathbf{f}_x^0, \dots, \mathbf{f}_x^\ell], \forall x \in \mathbf{P}$;
- 8 **end**
- 9 **return** $\{\mathbf{r}_x \leftarrow \mathbf{f}_x^L, \forall x \in \mathbf{P}\}$

where $\tilde{\phi}$, the grouped version of SLP ϕ , can widen its output to enhance its learning ability and maintain the original efficiency, and ψ , a normal SLP (shared over each centroid point x), is added to integrate the detached information in all groups. Both $\tilde{\phi}$ and ψ include a nonlinear activator.

To elaborate ePConv with filter grouping, let SLP $\tilde{\phi}$ (resp. SLP ψ) denote the SLP of $\tilde{\phi}$ (resp. ψ), and C_i (resp. C_o) denote the input (resp. output) channels of SLP $\tilde{\phi}$. N_g is the number of groups. Then, the parameter number of SLP $\tilde{\phi}$ before and after filter grouping is, $C_i \times C_o$ vs. $(C_i/N_g) \times (C_o/N_g) \times N_g = (C_i \times C_o)/N_g$. Here C_i and C_o are divisible by N_g and the few parameters in the bias term are ignored for clearness. In other words, using filter grouping, C_o can be increased by N_g times but with almost the same complexity. Besides, inspired by the bottleneck layer [9], we fix the output channels of SLP $\tilde{\phi}$ and SLP ψ as $C_o : k = 4 : 1$ (*i.e.*, $C_o = 4k$), to hold the original narrowness for DensePoint. Hence, with a small k , SLP ψ actually leads to only a little complexity of $4k \times k$, which can be easily remedied by a suitable N_g . The detailed forward pass procedure of DensePoint equipped with ePConv can be referred in Algorithm 1, where \ast indicates performing grouping operation.

DensePoint for point cloud processing. DensePoint applied in point cloud classification and per-point analysis (*e.g.*, segmentation) are illustrated in Fig. 3. In both tasks, DensePoint with ePConv is applied in each stage of the network to learn densely contextual representation, while PPool with original PConv is used to explicitly raise the semantic level and improve efficiency. For classification, the final global representation is learned by three P Pools and two DensePoints (11 layers in total, $L = 11$), followed by three fully connected (fc) layers as the classifier. For per-point analysis, four levels of representations learned by four P Pools and three DensePoints (17 layers in total, $L = 17$), are sequentially upsampled by feature propagation [34] to

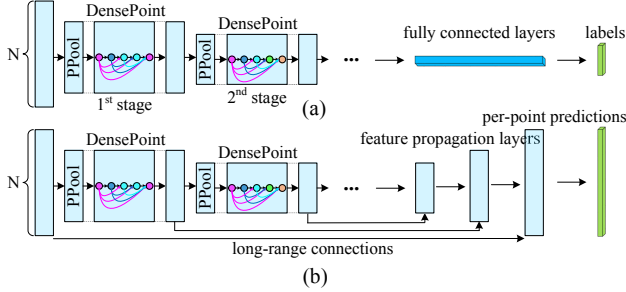


Figure 3. DensePoint applied in point cloud classification (a) and per-point analysis (b). PPool: pooling on point cloud (Sec 3.1). N is the number of points. The stage means several successive layers with the same number of points.

generate per-point predictions. All the networks can be trained in an end-to-end manner. The configuration details are included in the supplementary material.

Implementation details. PPool: the farthest points are picked from point cloud for uniform downsampling. Neighborhood: the spherical neighborhood is adopted; a fixed number of neighbors are randomly sampled in each neighborhood for batch processing (the centroid is reused if not enough), and they are normalized by subtracting the centroid. Group number N_g in $\tilde{\phi}$ (Eq. (4)): $N_g = 2$. Nonlinear activator: ReLU [29]. Dropout [42]: for model regularization, we apply dropout with 20% ratio on $f_{\mathcal{N}(x)}$ in Eq. (4) and dropout with 50% ratio on the first two fc layers in the classification network (Fig. 3(a)). Narrowness k : $k = 24$. Aggregation function ρ : symmetric function max pooling is employed. Batch normalization (BN) [17]: as done in image CNN, BN is used before each nonlinear activator for all layers. Note that only 3D coordinates (X, Y, Z) in \mathbb{R}^3 are used as the initial input features. Code is available¹.

4. Experiment

We conduct comprehensive experiments to validate the effectiveness of DensePoint. We first evaluate DensePoint for point cloud processing on challenging benchmarks across four tasks (Sec 4.1). We then provide detailed experiments to study DensePoint thoroughly (Sec 4.2).

4.1. DensePoint for Point Cloud Processing

Shape classification. We evaluate DensePoint on ModelNet40 and ModelNet10 classification benchmarks [53]. The former comprises 9843 training models and 2468 test models in 40 classes, while the latter consists of 3991 training models and 908 test models in 10 classes. The point cloud data is sampled from these models by [32]. For training, we uniformly sample 1024 points as the input. As in [21], we augment the input with random anisotropic scaling in range $[-0.66, 1.5]$ and translation in range $[-0.2, 0.2]$. For testing, similar to [32, 34], we apply voting with 10 tests using random scaling and then average the predictions.

¹<https://github.com/Yochengliu/DensePoint>

Table 1. Shape classification results (*overall accuracy, %*) on ModelNet40 (M40) and ModelNet10 (M10) benchmarks (pnt: point coordinates, nor: normal, “-”: unknown).

method	input	#points	M40	M10
Pointwise-CNN [12]	pnt	1k	86.1	-
Deep Sets [60]	pnt	1k	87.1	-
ECC [40]	pnt	1k	87.4	90.8
PointNet [32]	pnt	1k	89.2	-
SCN [55]	pnt	1k	90.0	-
Kd-Net(depth=10) [21]	pnt	1k	90.6	93.3
PointNet++ [34]	pnt	1k	90.7	-
MC-Conv [11]	pnt	1k	90.9	-
KCNet [39]	pnt	1k	91.0	94.4
MRTNet [4]	pnt	1k	91.2	-
Spec-GCN [49]	pnt	1k	91.5	-
DGCNN [52]	pnt	1k	92.2	-
PointCNN [26]	pnt	1k	92.2	-
PCNN [1]	pnt	1k	92.3	94.9
Ours	pnt	1k	93.2	96.6
SO-Net [24]	pnt	2k	90.9	94.1
Kd-Net(depth=15) [21]	pnt	32k	91.8	94.0
O-CNN [50]	pnt, nor	-	90.6	-
Spec-GCN [49]	pnt, nor	1k	91.8	-
PointNet++ [34]	pnt, nor	5k	91.9	-
SpiderCNN [56]	pnt, nor	5k	92.4	-
SO-Net [24]	pnt, nor	5k	93.4	95.7

Table 2. Shape retrieval results (mAP, %) on ModelNet40 (M40) and ModelNet10 (M10) benchmarks (“-”: unknown).

input	method	#points/views	M40	M10
Points	PointNet [10]	1k	70.5	-
	DGCNN [52]	1k	85.3	-
	PointCNN [26]	1k	83.8	-
	Ours	1k	88.5	93.2
Images	GVCNN [3]	12	85.7	-
	Triplet-center [10]	12	88.0	-
	PANORAMA-ENN [38]	-	86.3	93.3
	SeqViews [7]	12	89.1	89.5

The quantitative comparisons with the state-of-the-art point-based methods are summarized in Table 1. Our DensePoint outperforms all the point-input methods. Specifically, it reduces the error rate of PointNet++ by 26.9% on ModelNet40, and also surpasses its advanced version that applies additional normal data with very dense points (5k). Furthermore, even using only point as the input, DensePoint can also surpass the best additional-input method SO-Net [24] by 0.9% on ModelNet10. These results convincingly verify the effectiveness of DensePoint.

Shape retrieval. To further explore the recognition ability of DensePoint for the implicit shapes, we apply the global features, *i.e.*, the outputs of the penultimate fc layer in the classification network (Fig. 3(a)), for shape retrieval. We sort the most relevant shapes for each query from the test set by cosine distance, and report *mean Average Precision* (mAP). Except for point-based methods, we also compare with some advanced 2D image-based ones. The results are summarized in Table 2. As can be seen, DensePoint significantly outperforms PointNet by 18%. It is also com-

Table 3. Shape part segmentation results (%) on ShapeNet part benchmark (nor: normal, “-”: unknown).

method	input	class mIoU	instance mIoU	air plane	bag	cap	car	chair	ear phone	guitar	knife	lamp	laptop	motor bike	mug	pistol	rocket	skate board	table
Kd-Net [21]	4k	77.4	82.3	80.1	74.6	74.3	70.3	88.6	73.5	90.2	87.2	81.0	94.9	57.4	86.7	78.1	51.8	69.9	80.3
PointNet [32]	2k	80.4	83.7	83.4	78.7	82.5	74.9	89.6	73.0	91.5	85.9	80.8	95.3	65.2	93.0	81.2	57.9	72.8	80.6
SCN [55]	1k	81.8	84.6	83.8	80.8	83.5	79.3	90.5	69.8	91.7	86.5	82.9	96.0	69.2	93.8	82.5	62.9	74.4	80.8
SPLATNet [43]	-	82.0	84.6	81.9	83.9	88.6	79.5	90.1	73.5	91.3	84.7	84.5	96.3	69.7	95.0	81.7	59.2	70.4	81.3
KCNet [39]	2k	82.2	84.7	82.8	81.5	86.4	77.6	90.3	76.8	91.0	87.2	84.5	95.5	69.2	94.4	81.6	60.1	75.2	81.3
RS-Net [15]	-	81.4	84.9	82.7	86.4	84.1	78.2	90.4	69.3	91.4	87.0	83.5	95.4	66.0	92.6	81.8	56.1	75.8	82.2
DGCNN [52]	2k	82.3	85.1	84.2	83.7	84.4	77.1	90.9	78.5	91.5	87.3	82.9	96.0	67.8	93.3	82.6	59.7	75.5	82.0
PCNN [1]	2k	81.8	85.1	82.4	80.1	85.5	79.5	90.8	73.2	91.3	86.0	85.0	95.7	73.2	94.8	83.3	51.0	75.0	81.8
Ours	2k	84.2	86.4	84.0	85.4	90.0	79.2	91.1	81.6	91.5	87.5	84.7	95.9	74.3	94.6	82.9	64.6	76.8	83.7
SO-Net [24]	-,nor	80.8	84.6	81.9	83.5	84.8	78.1	90.8	72.2	90.1	83.6	82.3	95.2	69.3	94.2	80.0	51.6	72.1	82.6
SyncCNN [58]	mesh	82.0	84.7	81.6	81.7	81.9	75.2	90.2	74.9	93.0	86.1	84.7	95.6	66.7	92.7	81.6	60.6	82.9	82.1
PointNet++ [34]	2k,nor	81.9	85.1	82.4	79.0	87.7	77.3	90.8	71.8	91.0	85.9	83.7	95.3	71.6	94.1	81.3	58.7	76.4	82.6
SpiderCNN [56]	2k,nor	82.4	85.3	83.5	81.0	87.2	77.5	90.7	76.8	91.1	87.3	83.3	95.8	70.2	93.5	82.7	59.7	75.8	82.8

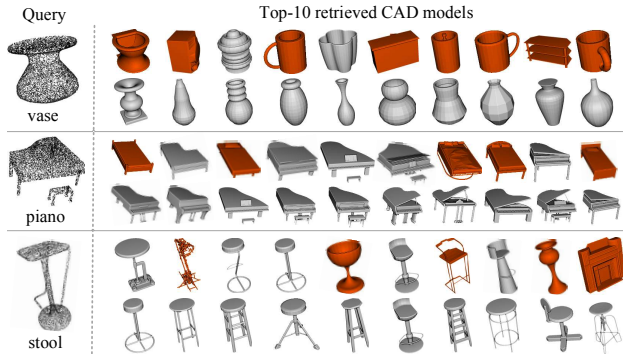


Figure 4. Retrieval examples on ModelNet40. Top-10 matches are shown for each query, with the 1st line for PointNet [32] and the 2nd line for our DensePoint. The mistakes are highlighted in red.

parable with those image-based methods (even the ensemble one [38]), which greatly benefit from image CNN and pre-training with large-scale datasets (e.g., ImageNet [37]). Fig. 4 shows some retrieval examples.

Shape part segmentation. Part segmentation is a challenging task for fine-grained shape recognition. Here we evaluate DensePoint on ShapeNet part benchmark [57]. It contains 16881 shapes with 16 categories, and is labeled in 50 parts in total, where each shape has 2~5 parts. We follow the data split in [32], and similarly, we also randomly pick 2048 points as the input and concatenate the one-hot encoding of the object label to the last feature layer of the segmentation network in Fig. 3(b). In testing, we also apply voting with ten tests using random scaling. Except for standard IoU (*Inter-over-Union*) score for each category, two types of mean IoU (mIoU) that are averaged across all classes and all instances respectively, are also reported.

Table 3 summarizes the quantitative comparisons with the state-of-the-art methods, where DensePoint achieves the best performance. Furthermore, it significantly surpasses the second best point-input methods, i.e., DGCNN [52], with 1.9 \uparrow in class mIoU and 1.3 \uparrow in instance mIoU respectively. Noticeably, it also sets new state of the arts over the point-based methods in eight categories. These improvements demonstrate the robustness of DensePoint to diverse

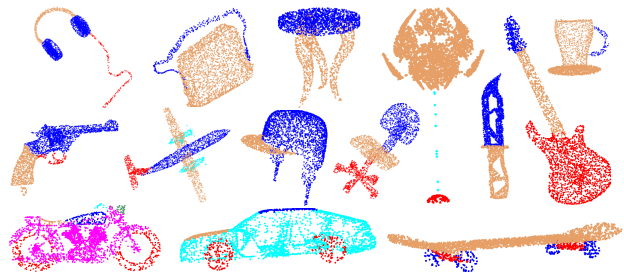


Figure 5. Segmentation examples on ShapeNet part benchmark.

Table 4. Normal estimation error on ModelNet40 benchmark.

dataset	method	#points	error
ModelNet40	PointNet [1]	1k	0.47
	PointNet++ [1]	1k	0.29
	PCNN [1]	1k	0.19
	MC-Conv [11]	1k	0.16
	Ours	1k	0.149

shapes. Some segmentation examples are shown in Fig. 5. **Normal estimation.** Normal estimation in point cloud is a crucial step for numerous applications, from surface reconstruction and scene understanding to rendering. Here, we regard normal estimation as a supervised regression task, and implement it by deploying DensePoint with the segmentation network in Fig. 3(b). The cosine-loss between the normalized output and the normal ground truth is employed for training. We evaluate DensePoint on ModelNet40 benchmark for this task, where 1024 points are uniformly sampled as the input.

The quantitative comparisons of the estimation error are summarized in Table 4, where DensePoint outperforms other advanced methods. Moreover, it significantly reduces the error of PointNet++ by 48.6%. Fig. 6 shows some normal prediction examples. As can be seen, DensePoint with densely contextual semantics can obtain more decent normal predictions, while PointNet and PointNet++ present a lot of deviations above 90 $^\circ$ from the ground truth. However, in this task, DensePoint can not process some intricate shapes well, e.g., curtains and plants.

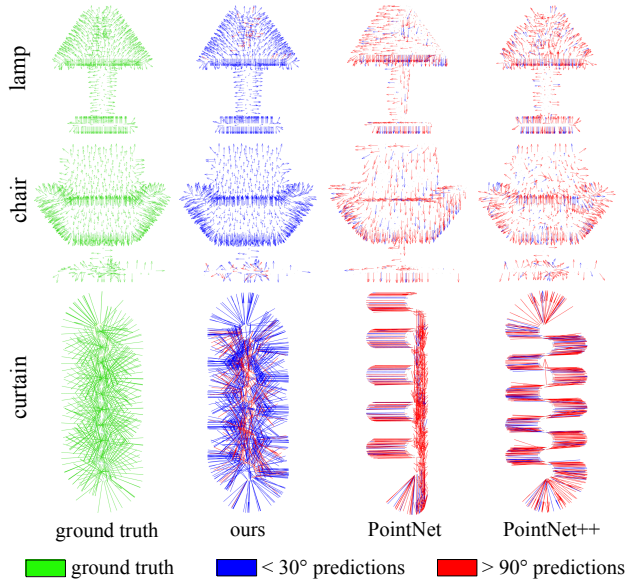


Figure 6. Normal estimation examples on ModelNet40 benchmark. For visual clarity, we only show the predictions with the angle less than 30° in blue, and the angle greater than 90° in red between the ground truth.

4.2. DensePoint Analysis

In this section, we first perform a detailed ablation study for DensePoint. Then, we discuss the group number N_g in ePConv (Eq. (4)), the network narrowness k , the aggregation function ρ and the network stage to apply DensePoint, respectively. Finally, we analyze the robustness of DensePoint on sampling density and random noise, and investigate the model complexity. All the experiments are conducted on ModelNet40 [53] dataset.

Ablation study. The results of ablation study are summarized in Table 5. We set two baselines: model A and model \bar{A} . Model A is set as a classic hierarchical version (the upper part of Fig. 2, *i.e.*, layer-by-layer connections without contextual learning by DensePoint) of the classification network with the same number of layers, and each layer is configured with the same width. Model \bar{A} directly concatenates all layers in each stage of model A as the final output of that stage. Both of them are equipped with PConv in Eq. (1).

The baseline model A gets a low classification accuracy of 88.6%, and it increases by only 0.5 percent with direct concatenation (model \bar{A}). However, with the densely contextual semantics of DensePoint, the accuracy raises significantly by 2.5 percent (91.1%, model B). This convincingly verifies its effectiveness. Then, when using ePConv to enhance the expressive power of each layer in DensePoint, the accuracy can be further improved to 92.5% (model C). Noticeably, the dropout on $\mathbf{f}_{\mathcal{N}(x)}$ in Eq. (4) can bring a boost of 0.3 percent (model D). The data augmentation technique can result in an accuracy variation of 0.7 percent (model E).

Table 5. Ablation study of DensePoint (%) (DA: data augmentation, DP: DensePoint, DO: dropout on $\mathbf{f}_{\mathcal{N}(x)}$ in Eq. (4)).

model	#points	DA	DP	ePConv	DO	vote	acc.
A	1k	✓					88.6
\bar{A}	1k	✓					89.1
B	1k	✓	✓				91.1
C	1k	✓	✓	✓			92.5
D	1k	✓	✓	✓	✓		92.8
E	1k	✓	✓	✓	✓		92.1
F	1k	✓	✓	✓	✓	✓	93.2
G	2k	✓	✓	✓	✓	✓	93.2

Table 6. The impact of the group number N_g on network parameters, FLOPs and performance ($k = 24$).

group number N_g	#params	#FLOPs/sample	acc. (%)
1	0.73M	1030M	92.7
2	0.67M	651M	93.2
4	0.62M	457M	92.2
6	0.61M	394M	92.3
12	0.60M	331M	92.1

Finally, by voting strategy, our final model F can achieve an impressive accuracy of 93.2%. In addition, we also investigate the number of input points by increasing it to 2k, yet obtaining no gain (model G). Maybe the model needs to be modified to adapt for more input points.

Group number N_g in ePConv (Eq. (4)). The filter grouping can greatly reduce the model complexity, whilst leading to a model regularization by rarefying the filter relationships [16]. Table 6 summarizes the impact of N_g on model parameters, model FLOPs (floating point operations/sample) and classification accuracy. As can be seen, the model parameters are very few (0.73M), even though the filter grouping is not performed. This is due to the narrow design ($k = 24$) of each layer in DensePoint and few parameters in the generalized convolution operator, ePConv. Eventually, with $N_g = 2$, DensePoint can achieve the best result of 93.2% with acceptable model complexity.

Network narrowness k . Table 7 summarizes the comparisons of different k . One can see that a very small DensePoint, *i.e.*, $k = 12$, can even obtain an impressive accuracy of 92.1%. This further verifies the powerfulness of the densely contextual semantics acquired by DensePoint on shape identification. Note that a large k is usually unnecessary for DensePoint, as it will greatly raise the model complexity but not bring any gains.

Aggregation function ρ . We experiment with three symmetric functions, *i.e.*, sum, average pooling and max pooling, whose results are 91.0%, 91.3% and 93.2%, respectively. The max pooling performs best, probably because it can select the biggest feature response to keep the most expressive representation and remove redundant information.

Network stage to apply DensePoint. To investigate the

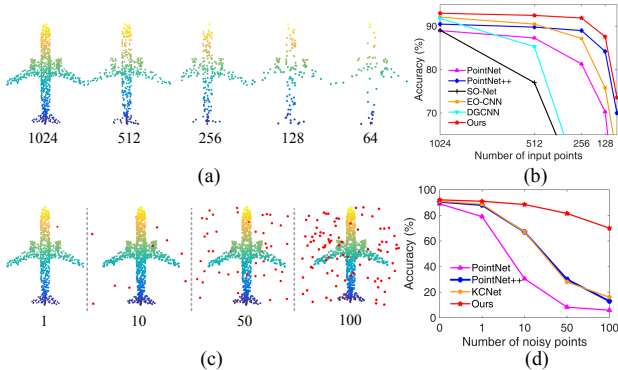


Figure 7. (a) Point cloud with different sampling densities. (b) Results of testing with sparser points. (c) Point cloud with some points being replaced with random noise (highlighted in red). (d) Results of testing with noisy points.

Table 7. The comparisons of different narrowness k ($N_g = 2$).

narrowness k	#params	#FLOPs/sample	acc. (%)
12	0.56M	294M	92.1
24	0.67M	651M	93.2
36	0.76M	957M	92.9
48	0.88M	1310M	92.7

Table 8. The comparisons of DensePoint applied in different stages of the classification network (Fig. 3(a)).

model	1 st stage	2 nd stage	acc. (%)
\tilde{A}			90.5
B	✓		91.8
C		✓	92.3
D	✓	✓	93.2

impact of contextual semantics at different levels on shape recognition, we also apply DensePoint with ePConv in different stages of the classification network (Fig. 3(a)). The results are summarized in Table 8. The baseline (model \tilde{A}) is set as the same as the model A in Table 5 but equipped with ePConv for a fair comparison. One can see that DensePoint applied in the 1st stage (model B) and the 2nd stage (model C) can both bring a considerable boost, while the latter performs better. This indicates the higher-level contextual semantics in the 2nd stage can result in a more powerful representation for shape recognition. Finally, with DensePoint in each stage for sufficiently contextual semantic information, the best result of 93.2% can be reached.

Robustness analysis. The robustness of DensePoint on sampling density and random noise are shown in Fig. 7. For the former, we use sparser points of 1024, 512, 256, 128 and 64, as the input to a model trained with 1024 points. Random input dropout is applied during training, for fair comparisons with PointNet [32], PointNet++ [34], SO-Net [24], PCNN [1] and DGCNN [52]. Fig. 7(b) shows that our model and PointNet++ perform better in this testing. Nevertheless, our model can obtain higher accuracy than Point-

Table 9. The comparisons of model complexity (“-”: unknown).

method	#params	#FLOPs/sample
PointNet [32]	3.50M	440M
PointNet++ [26]	1.48M	1684M
DGCNN [26]	1.84M	2767M
SpecGCN [26]	2.05M	1112M
KCNet [39]	0.90M	-
PCNN [26]	8.20M	294M
PointCNN [26]	0.60M	1581M
Ours ($k = 24, L = 11$)	0.67M	651M
Ours ($k = 24, L = 6$)	0.53M	148M

Net++ at all densities. This indicates the densely contextual semantics of DensePoint, is much more effective than the traditional multi-scale information of PointNet++.

For the latter, as in KCNet [39], we replace a certain number of randomly picked points with uniform noise ranging $[-1.0, 1.0]$ during testing. The comparisons with PointNet, PointNet++ and KCNet are shown in Fig. 7(d). Note that for this testing, our model is trained without any data augmentations to avoid confusion. As can be seen, our model is quite robust on random noise, while the others are vulnerable. This demonstrates the powerfulness of densely contextual semantics in DensePoint.

Model complexity. The comparisons of model complexity with the state of the arts are summarized in Table 9. As can be seen, our model is quite competitive and it can be the most efficient one with the network depth $L = 6$ (accuracy 92.1%). This shows its great potential for real-time applications, e.g., scene parsing in autonomous driving.

Discussion of limitations. (1) The density of local point clouds is not considered, which could lead to less effectiveness in greatly non-uniform distribution; (2) The importance of each level of context is not evaluated, which could lead to the difficulty in identifying very alike shapes.

5. Conclusion

In this work, DensePoint, a general architecture to learn densely contextual representation for efficient point cloud processing, has been proposed. DensePoint extends regular grid CNN to irregular point configuration by an efficient generalized convolution operator. Based on this operator, DensePoint develops a deep hierarchy and progressively aggregate multi-level and multi-scale semantics from it. As a consequence, DensePoint can acquire sufficiently contextual information along with rich semantics in an organic manner, making it highly effective for implicit shape identification. Extensive experiments on challenging benchmarks across four tasks, as well as thorough model analysis, have demonstrated that DensePoint achieves the state of the arts. In addition, DensePoint shows quite good robustness against noisy points, which could provide a promising direction for robust point cloud representation learning.

References

- [1] Matan Atzmon, Haggai Maron, and Yaron Lipman. Point convolutional neural networks by extension operators. In *SIGGRAPH*, pages 1–14, 2018. [2](#), [5](#), [6](#), [8](#)
- [2] Serge J. Belongie, Jitendra Malik, and Jan Puzicha. Shape matching and object recognition using shape contexts. *IEEE Trans. Pattern Anal. Mach. Intell.*, 24(4):509–522, 2002. [2](#)
- [3] Yifan Feng, Zizhao Zhang, Xibin Zhao, Rongrong Ji, and Yue Gao. GVCNN: Group-view convolutional neural networks for 3D shape recognition. In *CVPR*, pages 264–272, 2018. [1](#), [2](#), [5](#)
- [4] Matheus Gadelha, Rui Wang, and Subhransu Maji. Multiresolution tree networks for 3D point cloud processing. In *ECCV*, pages 105–122, 2018. [1](#), [2](#), [5](#)
- [5] Paul Guerrero, Yanir Kleiman, Maks Ovsjanikov, and Niloy J. Mitra. PCPNet: Learning local shape properties from raw point clouds. *Comput. Graph. Forum*, 37(2):75–85, 2018. [2](#)
- [6] Haiyun Guo, Jinqiao Wang, Yue Gao, Jianqiang Li, and Hanqing Lu. Multi-view 3D object retrieval with deep embedding network. *IEEE Trans. Image Processing*, 25(12):5526–5537, 2016. [1](#), [2](#)
- [7] Zhizhong Han, Mingyang Shang, Zhenbao Liu, Chi-Man Vong, Yu-Shen Liu, Matthias Zwicker, Junwei Han, and C. L. Philip Chen. SeqViews2SeqLabels: Learning 3D global features via aggregating sequential views by RNN with attention. *IEEE Trans. Image Processing*, 28(2):658–672, 2019. [2](#), [5](#)
- [8] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification. In *ICCV*, pages 1026–1034, 2015. [12](#)
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. [4](#)
- [10] Xinwei He, Yang Zhou, Zhichao Zhou, Song Bai, and Xiang Bai. Triplet-Center loss for multi-view 3D object retrieval. In *CVPR*, pages 1945–1954, 2018. [5](#)
- [11] Pedro Hermosilla, Tobias Ritschel, Pere-Pau Vázquez, Alvar Vinacua, and Timo Ropinski. Monte carlo convolution for learning on non-uniformly sampled point clouds. *ACM Trans. Graph.*, 37(6):235:1–235:12, 2018. [5](#), [6](#)
- [12] Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. Point-wise convolutional neural networks. In *CVPR*, pages 974–993, 2018. [5](#)
- [13] Gao Huang, Zhuang Liu, Laurens van der Maaten, and Kilian Q. Weinberger. Densely connected convolutional networks. In *CVPR*, pages 2261–2269, 2017. [2](#), [3](#), [4](#)
- [14] Haibin Huang, Evangelos Kalogerakis, Siddhartha Chaudhuri, Duygu Ceylan, Vladimir G. Kim, and Ersin Yumer. Learning local shape descriptors from part correspondences with multiview convolutional networks. *ACM Trans. Graph.*, 37(1):6:1–6:14, 2018. [2](#)
- [15] Qiangui Huang, Weiyue Wang, and Ulrich Neumann. Recurrent slice networks for 3D segmentation of point clouds. In *CVPR*, pages 2626–2635, 2018. [6](#)
- [16] Yani Ioannou, Duncan P. Robertson, Roberto Cipolla, and Antonio Criminisi. Deep roots: Improving CNN efficiency with hierarchical filter groups. In *CVPR*, pages 5977–5986, 2017. [7](#)
- [17] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, pages 448–456, 2015. [5](#)
- [18] Varun Jampani, Martin Kiefel, and Peter V. Gehler. Learning sparse high dimensional filters: Image filtering, dense crfs and bilateral neural networks. In *CVPR*, pages 4452–4461, 2016. [2](#)
- [19] Mingyang Jiang, Yiran Wu, and Cewu Lu. PointSIFT: A SIFT-like network module for 3D point cloud semantic segmentation. *arXiv preprint arXiv:1807.00652*, 2018. [1](#)
- [20] David Inkyu Kim and Gaurav S. Sukhatme. Semantic labeling of 3D point clouds with object affordance for robot manipulation. In *ICRA*, pages 5578–5584, 2014. [1](#)
- [21] Roman Klokov and Victor S. Lempitsky. Escape from cells: Deep Kd-Networks for the recognition of 3D point cloud models. In *ICCV*, pages 863–872, 2017. [2](#), [5](#), [6](#)
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In *NeurIPS*, pages 1106–1114, 2012. [1](#), [4](#)
- [23] Loic Landrieu and Martin Simonovsky. Large-scale point cloud semantic segmentation with superpoint graphs. In *CVPR*, pages 4558–4567, 2018. [2](#)
- [24] Jiaxin Li, Ben M. Chen, and Gim Hee Lee. SO-Net: Self-organizing network for point cloud analysis. In *CVPR*, pages 9397–9406, 2018. [5](#), [6](#), [8](#)
- [25] Ruoyu Li, Sheng Wang, Feiyun Zhu, and Junzhou Huang. Adaptive graph convolutional neural networks. In *AAAI*, pages 3546–3553, 2018. [2](#)
- [26] Yangyan Li, Rui Bu, Mingchao Sun, and Baoquan Chen. PointCNN: Convolution on X-transformed points. In *NeurIPS*, pages 828–838, 2018. [5](#), [8](#)
- [27] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *CVPR*, pages 8895–8904, 2019. [1](#)
- [28] Daniel Maturana and Sebastian Scherer. VoxNet: A 3D convolutional neural network for real-time object recognition. In *IROS*, pages 922–928, 2015. [1](#), [2](#)
- [29] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, pages 807–814, 2010. [5](#)
- [30] Aude Oliva and Antonio Torralba. The role of context in object recognition. *Trends in Cognitive Sciences*, 11(12):520–527, 2007. [1](#)
- [31] Charles R. Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J. Guibas. Frustum PointNets for 3D object detection from RGB-D data. In *CVPR*, pages 918–927, 2018. [1](#)
- [32] Charles R. Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. PointNet: Deep learning on point sets for 3D classification and segmentation. In *CVPR*, pages 77–85, 2016. [1](#), [2](#), [5](#), [6](#), [8](#)
- [33] Charles Ruizhongtai Qi, Hao Su, Matthias Nießner, Angela Dai, Mengyuan Yan, and Leonidas J. Guibas. Volumetric and multi-view CNNs for object classification on 3D data. In *CVPR*, pages 5648–5656, 2016. [2](#)
- [34] Charles R. Qi, Li Yi, Hao su, and Leonidas J. Guibas. PointNet++: Deep hierarchical feature learning on point sets in a

- metric space. In *NeurIPS*, pages 5099–5108, 2017. 1, 2, 3, 4, 5, 6, 8
- [35] Siamak Ravanbakhsh, Jeff Schneider, and Barnabas Poczos. Deep learning with sets and point clouds. In *ICLR*, pages 1–12, 2017. 1
- [36] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. OctNet: Learning deep 3D representations at high resolutions. In *CVPR*, pages 6620–6629, 2017. 2
- [37] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael S. Bernstein, Alexander C. Berg, and Fei-Fei Li. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 6
- [38] Konstantinos Sfikas, Ioannis Pratikakis, and Theoharis Theoharis. Ensemble of PANORAMA-based convolutional neural networks for 3D model classification and retrieval. *Computers & Graphics*, 71:208–218, 2018. 5, 6
- [39] Yiru Shen, Chen Feng, Yaoqing Yang, and Dong Tian. Mining point cloud local structures by kernel correlation and graph pooling. In *CVPR*, pages 4548–4557, 2018. 1, 2, 5, 6, 8
- [40] Martin Simonovsky and Nikos Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *CVPR*, pages 29–38, 2017. 5
- [41] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, pages 1–14, 2015. 1
- [42] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014. 5
- [43] Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. SPLATNet: Sparse lattice networks for point cloud processing. In *CVPR*, pages 2530–2539, 2018. 2, 6
- [44] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik G. Learned-Miller. Multi-view convolutional neural networks for 3D shape recognition. In *ICCV*, pages 945–953, 2015. 1, 2
- [45] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3D outputs. In *ICCV*, pages 2107–2115, 2017. 2
- [46] Gusi Te, Wei Hu, Amin Zheng, and Zongming Guo. RGCNN: Regularized graph CNN for point cloud segmentation. In *MM*, pages 746–754, 2018. 2
- [47] Du Tran, Lubomir D. Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3D convolutional networks. In *ICCV*, pages 4489–4497, 2015. 2
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 6000–6010, 2017. 2
- [49] Chu Wang, Babak Samari, and Kaleem Siddiqi. Local spectral graph convolution for point set feature learning. In *ECCV*, pages 1–16, 2018. 1, 2, 5
- [50] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-CNN: octree-based convolutional neural networks for 3D shape analysis. *ACM Trans. Graph.*, 36(4):72:1–72:11, 2017. 2, 5
- [51] Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong. Adaptive O-CNN: a patch-based deep representation of 3D shapes. *ACM Trans. Graph.*, 37(6):217:1–217:11, 2018. 2
- [52] 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 Trans. Graph.*, pages 1–13, 2019. 5, 6, 8
- [53] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D ShapeNets: A deep representation for volumetric shapes. In *CVPR*, pages 1912–1920, 2015. 1, 2, 5, 7
- [54] Jin Xie, Guoxian Dai, Fan Zhu, Edward K. Wong, and Yi Fang. DeepShape: Deep-learned shape descriptor for 3D shape retrieval. *IEEE Trans. Pattern Anal. Mach. Intell.*, 39(7):1335–1345, 2017. 2
- [55] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional ShapeContextNet for point cloud recognition. In *CVPR*, pages 4606–4615, 2018. 2, 5, 6
- [56] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. SpiderCNN: Deep learning on point sets with parameterized convolutional filters. In *ECCV*, pages 90–105, 2018. 5, 6
- [57] Li Yi, Vladimir G. Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas J. Guibas. A scalable active framework for region annotation in 3D shape collections. *ACM Trans. Graph.*, 35(6):210:1–210:12, 2016. 6
- [58] Li Yi, Hao Su, Xingwen Guo, and Leonidas J. Guibas. Sync-SpecCNN: Synchronized spectral CNN for 3D shape segmentation. In *CVPR*, pages 6584–6592, 2017. 6
- [59] Kangxue Yin, Hui Huang, Daniel Cohen-Or, and Hao (Richard) Zhang. P2P-NET: bidirectional point displacement net for shape transform. *ACM Trans. Graph.*, 37(4):152:1–152:13, 2018. 1
- [60] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabás Póczos, Ruslan R. Salakhutdinov, and Alexander J. Smola. Deep sets. In *NeurIPS*, pages 3394–3404, 2017. 5
- [61] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *ECCV*, pages 818–833, 2014. 1, 3

Supplementary Material

A. Outline

This supplementary material provides: (1) further investigations of the proposed DensePoint (Sec B); (2) more shape retrieval examples of DensePoint and some analysis (Sec C); (3) network configuration details (Sec D). (4) training details (Sec E)

B. Further Investigations

In this section, we provide further investigations of DensePoint on four aspects. Specifically, the discussion of neighborhood method is presented in Sec B.1. The effect of dropout on $f_{\mathcal{N}(x)}$ in Eq. (4) is analyzed in Sec B.2. The impact of network depth on classification performance is investigated in Sec B.3. The memory and runtime are summarized in Sec B.4. All the investigations are conducted on ModelNet40 dataset.

B.1. Neighborhood Method

In the main paper, the local convolutional neighborhood $\mathcal{N}(x)$ in Eq. (1) is set to be a spherical neighborhood, from which a fixed number of neighbors are randomly sampled for batch processing. We compare this strategy (Random-In-Sphere) with another typical one, *i.e.*, k-nearest neighbor (k-NN). For a fair comparison, the models with these two strategies are configured with the same settings. Table I summarizes the results.

As can be seen, the model with Random-In-Sphere performs better. We speculate that the model with k-NN will suffer from the distribution inhomogeneity of points. In this case, the contextual learning in DensePoint will be less effective, as the receptive fields will be confined to a local region with large density, which leads to ignoring those sparse points that are essential for recognizing the implicit shape. By contrast, Random-In-Sphere can have a better coverage of points even in the case of inhomogeneous distribution.

B.2. Dropout on $f_{\mathcal{N}(x)}$ in Eq. (4)

The dropout technique can force the whole network to behave as an ensemble of a lot of subsets and reduce the risk of model overfitting. To analyze its effect on DensePoint, we apply it with different ratios on $f_{\mathcal{N}(x)}$ in Eq. (4). The results are summarized in Table II. As can be seen, the best result of 93.2% can be achieved with a dropout ratio of 20%.

B.3. Network Depth

We further explore the impact of the network depth (fully connected layers are not included) on classification performance. The results are summarized in Table III. Surprisingly, a 6-layer network equipped with DensePoint can achieve an accuracy of 92.1% with only 0.53M params and 148M FLOPs/sample. This even outperforms PointNet++

Table I. The results (%) of two neighborhood strategies. The number of neighbors is equally set in each layer of the two models.

neighborhood method	acc.
k-NN	91.3
Random-In-Sphere	93.2

Table II. The results (%) of dropout with different ratios applied on $f_{\mathcal{N}(x)}$ in Eq. (4).

ratio (%)	0	10	20	30	40	50
acc.	92.9	92.8	93.2	93.0	92.8	92.5

Table III. The results (%) of different network depths (fully connected layers are not included).

#layers	#params	#FLOPs/sample	acc.
6	0.53M	148M	92.1
9	0.56M	510M	92.9
11	0.67M	651M	93.2
15	0.78M	779M	93.0
19	0.88M	1222M	92.7
23	1.03M	1416M	92.6

Table IV. Time and memory of classification network, where k is network narrowness, L is network depth. The statistics of all the models are summarized with batch size 16 on NVIDIA TITAN Xp, and time is the mean time of 1000 tests. The compared models are tested using their available official codes.

method	#points	Time (ms)		Memory (GB)	
		training	test	training	test
PointNet [31]	1024	55	22	1.318	0.469
PointNet++ [33]	1024	195	47	8.311	2.305
DGCNN [52]	1024	300	68	4.323	1.235
PointCNN [26]	1024	55	38	2.501	1.493
Ours ($k=24, L=11$)	1024	21	10	3.745	1.228
Ours ($k=24, L=6$)	1024	10	5	1.468	0.886

[33] (accuracy 90.7%, params 1.48M [26], FLOPs/sample 1684M [26]) by 15% in error rate, whilst being one order of magnitude faster in terms of FLOPs/sample. We also observe that it is unnecessary to develop a very deep network (*e.g.*, 23 layers) with DensePoint, as it increases complexity without bringing any gain. Eventually, the best result of 93.2% can be reached with acceptable complexity by an 11-layer network.

B.4. Memory and runtime

The memory and runtime of the proposed DensePoint are summarized in Table IV. As can be seen, the model ($L=11$) is competitive while another model ($L=6$) is the best one in terms of efficiency. Actually, the memory and training time issues in dense connection mode are greatly alleviated due to the shallow design of DensePoint and our highly-efficient implementation. Moreover, although extremely deep network could be unnecessary for 3D currently, in case of very deep DensePoint in the future, the technique of Shared Memory Allocations can be applied to achieve linear memory complexity.

C. Shape Retrieval

In this section, we show more shape retrieval examples in Fig. 8. As can be seen, compared with PointNet [31], our DensePoint obtains superior shape identification results. Specifically, PointNet is confused between the query “bottle” and the sample “vase” due to their similar shapes. Nevertheless, DensePoint with densely contextual semantics acquired can identify them accurately. We notice that DensePoint could also be confused for some very alike shapes, e.g., the query “bench” and the sample “tv_stand”. This could be improved by learning to weight multi-level contextual information instead of identically aggregating all levels of information. We leave it as future work.

D. Network Configuration Details

In this section, we present the configuration details of three networks on shape classification, shape part segmentation and normal estimation, respectively. For clearness, we describe the layer and corresponding setting format as follows:

PPool: [downsampling rate, neighborhood radius, #number of neighbors, SLP^ϕ (#input channels, #output channels)]. The global pooling is achieved by directly applying PConv to convolve all points.

ePConv: [neighborhood radius, #number of neighbors, SLP^ϕ (#input channels, #output channels, #group number), SLP^ψ (#input channels, #output channels), dropout ratio].

FP (feature propagation layer): $MLP(\#channels, \dots)$. Feature propagation layer [33] is used for transforming the features that are concatenated from current interpolated layer and long-range connected layer. We employ a multi-layer perceptron (MLP) to implement this transformation.

FC (fully connected layer): [(#input channels, #output channels), dropout ratio]. Note that the dropout technique is applied for all FC layers except for the last FC layer (used for prediction).

In addition, except for the last prediction layer, all layers (including the inside perceptrons) are followed with batch normalization and ReLU activator. The output shape is in the format of (#feature dimension, #number of points).

D.1. Shape Classification Network

The configuration details of shape classification network are presented in Table VI. The network has 14 layers in total, which comprises 3 Pools (the last one is global pooling layer) and 2 DensePoints (the 1st one has 3 layers while the 2nd one has 5 layers), followed by 3 FC layers.

D.2. Shape Part Segmentation Network

Table V summarizes the configuration details of shape part segmentation network. As it shows, the network has

23 layers in total, which comprises 4 Pools, 3 DensePoints (4 layers, 6 layers and 3 layers in 2nd stage, 3rd stage and 4th stage respectively) and 4 FP layers, followed by 2 FC layers. As in [31, 33], we concatenate the one-hot encoding (16-d) of the object label to the last feature layer.

D.3. Normal Estimation Network

The normal estimation network is presented in Table VII. It is almost the same as the segmentation network, except for three aspects: (1) the input becomes 1024-d and the one-hot encoding becomes 40-d for ModelNet40 dataset; (2) the settings of some layers are slightly changed to be consistent with the 1024-d input; (3) the final output becomes 3-d for normal prediction. As done in the segmentation network, we also concatenate the one-hot encoding (40-d) of the object label to the last feature layer.

E. Training Details

Our DensePoint is implemented using Pytorch. The Adam optimization algorithm is employed for training, with a mini-batch size of 32. The momentum for batch normalization starts with 0.9 and decays with a rate of 0.5 every 20 epochs. The learning rate begins with 0.001 and decays with a rate of 0.7 every 20 epochs. The weight is initialized using the techniques introduced by He *et al.* [8].

Table V. The configuration details of shape part segmentation network. “long-range” indicates the long-range connections (see Fig. 3(b) in the main paper). K is the number of classes.

stage	layer type	setting detail	output shape	long-range
-	Input	-	(3, 2048)	FP ₄
1	PPool	[1/2, 0.1, 32, (3, 64)]	(64, 1024)	FP ₃
	PPool	[1/4, 0.2, 64, (64, 128)]	(128, 256)	
2	ePConv	[0.3, 32, (128, 96, 2), (96, 24), 20%]	(24, 256)	
	ePConv	[0.3, 32, (152, 96, 2), (96, 24), 20%]	(24, 256)	
	ePConv	[0.3, 32, (176, 96, 2), (96, 24), 20%]	(24, 256)	
	ePConv	[0.3, 32, (200, 96, 2), (96, 24), 20%]	(24, 256)	
The output of DensePoint in 2 nd stage			(224, 256)	FP ₂
	PPool	[1/4, 0.3, 32, (224, 192)]	(192, 64)	
3	ePConv	[0.5, 16, (192, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (216, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (240, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (264, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (288, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (312, 96, 2), (96, 24), 20%]	(24, 64)	
The output of DensePoint in 3 rd stage			(336, 64)	FP ₁
	PPool	[1/4, 0.8, 32, (336, 360)]	(360, 16)	
4	ePConv	[0.8, 8, (360, 96, 2), (96, 24), 20%]	(24, 16)	
	ePConv	[0.8, 8, (384, 96, 2), (96, 24), 20%]	(24, 16)	
	ePConv	[0.8, 8, (408, 96, 2), (96, 24), 20%]	(24, 16)	
The output of DensePoint in 4 th stage			(432, 16)	
	FP ₁	(768, 512, 512)	(512, 64)	
	FP ₂	(736, 384, 384)	(384, 256)	
	FP ₃	(448, 256, 256)	(256, 1024)	
	FP ₄	(259, 128, 128)	(128, 2048)	
	FC	[(128+16 ¹ , 128), 50%]	(128, 2048)	
	FC	[(128, K), -] → softmax	(K , 2048)	

¹ This is the one-hot encoding of the object label on ShapeNet part dataset.

Table VI. The configuration details of shape classification network. K is the number of classes.

stage	layer type	setting detail	output shape
-	Input	-	(3, 1024)
	PPool	[1/2, 0.25, 64, (3, 96)]	(96, 512)
1	ePConv	[0.2, 32, (96, 96, 2), (96, 24), 20%]	(24, 512)
	ePConv	[0.2, 32, (120, 96, 2), (96, 24), 20%]	(24, 512)
	ePConv	[0.2, 32, (144, 96, 2), (96, 24), 20%]	(24, 512)
The output of DensePoint in 1 st stage			(168, 512)
	PPool	[1/4, 0.3, 64, (168, 144)]	(144, 128)
2	ePConv	[0.4, 16, (144, 96, 2), (96, 24), 20%]	(24, 128)
	ePConv	[0.4, 16, (168, 96, 2), (96, 24), 20%]	(24, 128)
	ePConv	[0.4, 16, (192, 96, 2), (96, 24), 20%]	(24, 128)
	ePConv	[0.4, 16, (216, 96, 2), (96, 24), 20%]	(24, 128)
	ePConv	[0.4, 16, (240, 96, 2), (96, 24), 20%]	(24, 128)
The output of DensePoint in 2 nd stage			(264, 128)
	PPool	[-, -, 128, (264, 512)]	(512,)
3	FC	[(512, 512), 50%]	(512,)
	FC	[(512, 256), 50%]	(256,)
	FC	[(256, K), -] → softmax	(K ,)



Figure 8. Retrieval examples on ModelNet40 dataset. Top-10 matches are shown for each query, with the 1st line for PointNet [31] and the 2nd line for our DensePoint. The mistakes are highlighted in red.

Table VII. The configuration details of normal estimation network. “long-range” indicates the long-range connections (see Fig. 3(b) in the main paper).

stage	layer type	setting detail	output shape	long-range
-	Input	-	(3, 1024)	FP ₄
1	PPool	[1, 0.2, 32, (3, 64)]	(64, 1024)	FP ₃
	PPool	[1/4, 0.2, 32, (64, 128)]	(128, 256)	
2	ePConv	[0.3, 32, (128, 96, 2), (96, 24), 20%]	(24, 256)	
	ePConv	[0.3, 32, (152, 96, 2), (96, 24), 20%]	(24, 256)	
	ePConv	[0.3, 32, (176, 96, 2), (96, 24), 20%]	(24, 256)	
	ePConv	[0.3, 32, (200, 96, 2), (96, 24), 20%]	(24, 256)	
	The output of DensePoint in 2 nd stage		(224, 256)	FP ₂
3	PPool	[1/4, 0.3, 32, (224, 192)]	(192, 64)	
	ePConv	[0.5, 16, (192, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (216, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (240, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (264, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (288, 96, 2), (96, 24), 20%]	(24, 64)	
	ePConv	[0.5, 16, (312, 96, 2), (96, 24), 20%]	(24, 64)	
	The output of DensePoint in 3 rd stage		(336, 64)	FP ₁
4	PPool	[1/4, 0.8, 32, (336, 360)]	(360, 16)	
	ePConv	[0.8, 8, (360, 96, 2), (96, 24), 20%]	(24, 16)	
	ePConv	[0.8, 8, (384, 96, 2), (96, 24), 20%]	(24, 16)	
	ePConv	[0.8, 8, (408, 96, 2), (96, 24), 20%]	(24, 16)	
	The output of DensePoint in 4 th stage		(432, 16)	
	FP ₁	(768, 512, 512)	(512, 64)	
	FP ₂	(736, 384, 384)	(384, 256)	
	FP ₃	(448, 256, 256)	(256, 1024)	
	FP ₄	(259, 128, 128)	(128, 1024)	
	FC	[(128+40 ² , 128), 50%]	(128, 1024)	
	FC	[(128, 3), -]	(3, 1024)	

² This is the one-hot encoding of the object label on ModelNet40 dataset.