

Research Article

A New Volumetric CNN for 3D Object Classification Based on Joint Multiscale Feature and Subvolume Supervised Learning Approaches

A. A. M. Muzahid ^{1,2}, Wanggen Wan,^{1,2} and Li Hou³

¹*School of Communication and Information Engineering, Shanghai University, Shanghai 200444, China*

²*Institute of Smart City, Shanghai University, Shanghai 200444, China*

³*School of Information Engineering, Huangshan University, Anhui 245041, China*

Correspondence should be addressed to A. A. M. Muzahid; muzahid@shu.edu.cn

Received 12 September 2019; Revised 2 May 2020; Accepted 7 May 2020; Published 28 May 2020

Academic Editor: Daniele Bibbo

Copyright © 2020 A. A. M. Muzahid et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

The advancement of low-cost RGB-D and LiDAR three-dimensional (3D) sensors has permitted the obtainment of the 3D model easier in real-time. However, making intricate 3D features is crucial for the advancement of 3D object classifications. The existing volumetric voxel-based CNN approaches have achieved remarkable progress, but they generate huge computational overhead that limits the extraction of global features at higher resolutions of 3D objects. In this paper, a low-cost 3D volumetric deep convolutional neural network is proposed for 3D object classification based on joint multiscale hierarchical and subvolume supervised learning strategies. Our proposed deep neural network inputs 3D data, which are preprocessed by implementing memory-efficient octree representation, and we propose to limit the full layer octree depth to a certain level based on the predefined input volume resolution for storing high-precision contour features. Multiscale features are concatenated from multilevel octree depths inside the network, aiming to adaptively generate high-level global features. The strategy of the subvolume supervision approach is to train the network on subparts of the 3D object in order to learn local features. Our framework has been evaluated with two publicly available 3D repositories. Experimental results demonstrate the effectiveness of our proposed method where the classification accuracy is improved in comparison to existing volumetric approaches, and the memory consumption ratio and run-time are significantly reduced.

1. Introduction

Three-dimensional (3D) objects have a prevalent significance within the areas of computer vision applications including human-machine interactions and autonomous vehicles to robotics [1]. Deep learning (DL) achieved impressive success in 2D fields [2–4] with various applications such as face recognition and image classification. In human vision, whatever we see in the real world is within 3D space, so 3D data can improve the performance of computer vision-based applications [5]. In the last few years, there have been several 3D databases published to the public [6–8]. These have opened the door for computer vision researchers to work with real-world objects, and DL-based 3D shape analysis research has become possible, including 3D

classification, segmentation, retrieval, and 3D reconstruction. However, unlike the regular sampled 2D images, 3D shapes are irregular triangle meshes or point clouds; it is a challenging task for CNN to extract distinctive features [9] that can characterize the shapes and parts of a 3D object. In this paper, the 3D object classification task is considered by employing volumetric deep convolutional neural network (CNN) using 3D CAD models. Figure 1 shows a general block diagram of the 3D object classification.

One of the earliest attempts of CNNs to recognize a 3D object was revealed using depth information of RGB-D images [10]. In recent years, a number of papers have been published using CNN-based approaches for 3D object classification tasks. Among them, 2D-based CNN approaches obtained popularity as existing 2D-based CNN

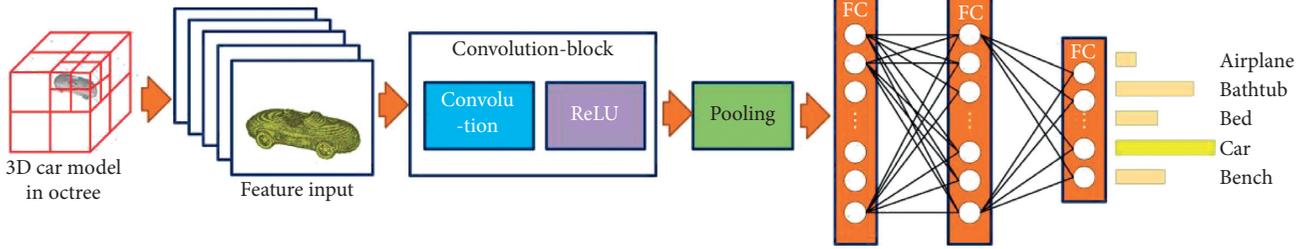


FIGURE 1: A basic block diagram of CNN-based 3D object classification.

frameworks could be used directly that required lower computational costs; besides, they achieved a more precise result when a 3D object was projected to 2D images [11–14], commonly known as multiview CNN. However, multiple projections of a 3D object into 2D grids discard intrinsic information, and these are not compatible with complex 3D vision tasks. To capture the complete geometric features of a 3D object, volumetric representation with periodic grid-style data offers very comprehensive information where 3D data can be sent directly to volumetric CNNs. Voxels are commonly used to exploit the 3D model in volumetric representation and directly fed to CNNs. 3DShapeNets [15] are the first volumetric CNNs proposed by Wu et al., and they have released a rich 3D CAD model repository of ModelNet datasets [6]. 3DShapeNets adopted Convolution Deep Belief Net (CDBN) from 2D DL to 3D distribution and were applied for three different applications, including 3D object classifications, next view predictions, and retrieval tasks. A similar approach was used in [16], which proposed a shallow volumetric CNN named VoxNet for real-time 3D object recognition tasks from three different domains. VoxNet achieved faster recognition capabilities and outperformed 3DShapeNets for an object recognition task on the ModelNet40 dataset. In the following, a few more approaches are proposed that use volumetric CNNs [17–20] and achieved state-of-the-art performance for 3D object classification tasks. In spite of the effectiveness of voxel-based CNNs, they generate a huge computational overhead due to the convolution operation of unnecessary bounding voxel data with many network parameters. In 2D CNNs, a 256×256 image resolution is normally considered as the input for DL-based image classification, but considering a 3D object with a 256^3 volume resolution is computationally prohibited. Although the computational power of GPU-based CNNs has significantly improved recently, in contrast, the training time and computational issues are the main constraints of voxel-based 3D volumetric data that limit the use of high-volume resolutions and going for deeper networks. The octree-based volumetric representation started gaining popularity by researchers because it reduced computational overhead [21–24]. However, octree representation performs better to preserve the fine details of the 3D object and the smoothness of the 3D object’s surface in comparison to voxel representation.

In this paper, we also consider the low-cost octree as a volumetric representation of the 3D object, which has a long history of different 3D data applications [25–29]. This octree representation is formed on a recursive decomposition of the

root voxels where the data tree divides the total 3D volume into cubes similar to quadtree structures [30]. Sparsely occupied cubes by the object are inserted into the data tensor where the CNN is performed, and the computational cost increases quadratically. We propose a new multiscale volumetric deep convolutional neural network (MS-VDCNN) based on joint residual and subvolume supervised learning approaches using an octree data tensor as the input of the network. In comparison to regular octree representation, we propose to reserve the full-layer octree depth to a certain minimum octree level, depending on the predefined voxel resolution. By this phenomenon, all volume features will be frozen until the system-defined minimum octree level is achieved. This reservation will help to address information loss while converting features from 1D octree to 3D volume space inside our network. A shallow multitask learning framework was adopted to reduce training errors and optimize the system shortly. This residual learning block concatenates multilevel conv. features to improve performance. In addition, the subvolume supervision was employed using layer slicing tactics where the object volume is divided into subvolume parts, and the network is trained on sliced data tensor to learn local features. This approach can be compared with our human vision tactics, as it has the ability to identify an object by observing some parts of it. Our key contributions of this paper are threefold: proposing a unified volumetric framework, effective octree structures, and an investigation of optimal training samples. These can be summarized as follows:

- (i) We propose a unified GPU-based MS-VDCNN for volumetric object classification. This proposed method aims to make full use of multilevel features at the higher resolution of input samples (3D objects) utilizing a residual hierarchical learning approach. A subvolume supervised learning tactic is applied to tackle the overfitting issue and improve the network performance.
- (ii) Our MS-VDCNN straightforwardly inputs the octree data tensor and generates feature maps based on joint residual and subvolume supervised learning methods. An optimal octree representation is formed by reserving full octants to a certain octree depth based on the predefined sample resolution. In this way, the octree encoder stores high-precision global features at the beginning of the octree partition that helps to improve the performance of the network. This approach may seem to increase the memory consumption ratio slightly in comparison to regular octree

representation, but it still consumes less memory than full voxel methods (Figure 2).

- (iii) The influence of input volume resolutions and multiorientation effects is investigated by employing extensive experiments on the ModelNet40 dataset. Based on the experimental results, the optimal volume resolution and number of views are defined to further improve the network performance in terms of classification accuracy and loss estimation. Our proposed MS-VDCNN achieves higher classification accuracy compared with other single volumetric CNNs.

The rest of this paper is organized as follows: Section 2 presents earlier CNN based related works for 3D object classification. In Section 3, our proposed methods including volumetric point-octree representation and network architecture are discussed. Section 4 presents the experimental results. Finally, Section 5 summarizes our conclusion and future works.

2. Related Work

From the motivational results achieved by the DL methods [2, 3, 31] in the 2D field, recently, DL has also attracted the attention of the 3D computer vision community in order to learn the complex structures of 3D data. A set of 3D deep neural networks has been published, during the last few years, for several 3D vision applications, including 3D object classification, shape retrieval, and part segmentation. However, the performance of DL networks strongly depends on the representation of 3D data, the design of CNN frameworks, and generalization of network parameters. In this section, some existing DL models, including shape descriptors, volumetric, 2D projections, and multiview and point cloud approaches, will be reviewed for 3D object classification.

3D shape descriptors provide some key properties of the 3D shape. Based on the type of feature extraction, a descriptor can be categorized into global and local descriptors. The whole 3D shape information can be extracted by global descriptor; by contrast, local descriptor represents the local surface with low-level features. The global descriptor cannot preserve the local details. Therefore, the local descriptor is commonly used to train a CNN framework to generate high-level features to learn the hierarchical discriminative features of the 3D shape for 3D computer vision applications [32, 33]. Liu et al. [34] proposed a local 3D shape descriptor to extract low-level visual features encoded into bag-of-words (BOW) from 200 multiview of a 3D shape. Furthermore, deep belief networks (DBNs) input the BOW paradigm of each shape to learn the high-level semantic features for classification and retrieval tasks. Experimental results demonstrated that the trained model achieved better results when compared with the classical BOW features. Han et al. [35] proposed deep spatiality (DS) as an unsupervised learning framework to learn both global and local features using deep neural networks. The spatial context extractor and the deep context learner are two major components of this framework. The

spatial relation in the local region was encoded by the spatial context extractor, and deep context learner was trained by both the global and local features using the coupled Softmax layer proposed by them. The low-level features of 3D object faces are used for classification and retrieval of tasks in a semisupervised or unsupervised manner. Recently, some more advanced approaches are published including adversarial attack for volumetric data [36] and energy-based model for generating volumetric shape patterns [37]. However, descriptor-based CNNs present 3D data with a form of abstraction which might not be effective for supervised learning. That is why descriptor features are mostly used in an unsupervised or semisupervised manner [38–41].

3D data projection is another representation of 3D data where some key properties of the original 3D raw data are projected into 2D space. Shi et al. [13] proposed DeepPano for 3D object recognition and retrieval tasks. A cylindrical 2D projection around the principle direction of the 3D object was made to extract 2D panoramic views. Cao et al. [42] proposed two complementary projections on a spherical domain, and this produced cylindrical patches. The first projection stores depth variations, and contour-information are encoded on different angles by the second projection. The proposed 2D CNN inputs a set of cylindrical patches as features of a 3D object. The trained model was applied for 3D object classifications. The major advantage of this projection method is to use 2D DL directly for 3D applications [12]. The performance of this method highly depends on the type of the projection method being used. However, this kind of representation is not optimal for 3D vision tasks because of information loss in the projection [43].

Multiview CNNs are another popular method for 3D vision applications. A 3D object is represented as a set of 2D images which allows for learning multiple features by employing existing 2D DL architectures. For computer vision applications, the surface information of a 3D object provides very high-level features that can be generated with multiple 2D views of a 3D object. Su et al. [11] proposed a multiview convolution neural network (MVCNN) for 3D object recognition. A pretrained model from ImageNet [44] was fine-tuned on the ModelNet40 dataset using 2D views (12 or 80 rendered images) of every 3D object as input features. The classification results on the ModelNet40 dataset by MVCNN achieved a state-of-the-art performance. MVCNN was further improved by Johns et al. [45] who proposed a multiview object recognition over arbitrary camera trajectories. An image sequence was decomposed into a set of pair, and each pair was classified independently. Ma et al. proposed a novel multiview model [46] that extracts the correlation between multiviews by combining CNNs with long short-term memory (LSTM). Ha et al. introduced 3D2SeqViews [47] that generate global features by aggregating sequential views and perform better than other correlation-based methods by Ma et al. [46]. In general, multiview-based CNNs give better performances when compared with other 3D representation methods. Therefore, the number of views is still an open issue for multiview representation where insufficient views may result

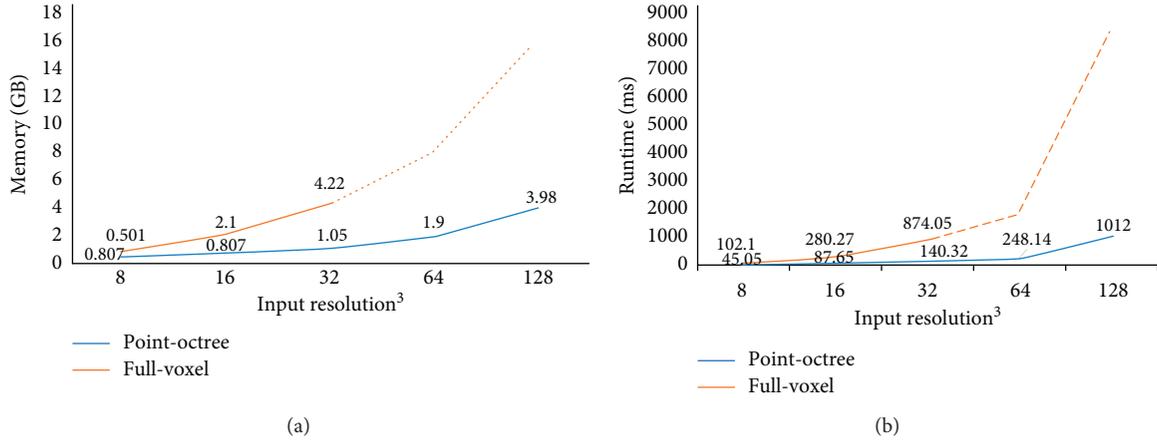


FIGURE 2: Comparison results between our point-octree and full voxel methods on ModelNet40 dataset. The total GPU memory was 5 GB and the batch size was 32: (a) memory consumption ratio in GB and (b) average runtime in milliseconds.

in underfitting issues, and too many views may increase the computational cost [48].

A Point cloud represents the 3D object geometry as a set of unstructured 3D points. However, point cloud CNNs are becoming popular as they can input point cloud data directly, and a lot of work has been done in multiple computer vision tasks (e.g., object recognition and 3D reconstruction) [49–51]. In this domain, PointNet [50], the pioneer deep net framework, was trained on unordered point sets in 3D environments to perform classification and segmentation tasks. The efficient PointNet achieves a state-of-the-art performance. Qi et al. proposed PointNet++ [52] where PointNet was recursively applied on a nested partitioning of the input point set. The PointNet++ framework takes into account the partitioning of the point set, and it extracts local features from PointNet. PointNet++ performs better than PointNet and obtained the benchmarks of 3D point clouds. In the following, a few more point cloud CNNs have been published, and they have improved 3D classification performance, e.g., Kd-network [49] presents the Kd-tree structures of point cloud data; PVNet [53] learns 3D geometric features from the point clouds and rendered images; relation-shape CNN (RS-CNN) introduces contextual shape-aware learning [54]; 3D Capsule [55] presents a new module, namely, ComposeCap, to map spatially relevant feature for point cloud classifications. However, point cloud CNNs face some difficulties because of the spatial irregular and permutation invariant properties [1]. Additional processing might be required to generate high-quality surface features due to the lack of connectivity information.

Volumetric representations give the full geometry of a 3D object. 3DShapeNets [15] is a very early attempt at volumetric CNNs for 3D shape analysis using 3D CAD models, and the released ModelNet datasets [6] has become a benchmark for 3D object classification. CDBN was adopted from 2D DL that inputs a binary voxel tensor with a volume resolution of $30 \times 30 \times 30$. However, the structure and geometric properties of the 3D object can vary on

different representations. To exploit these properties, VoxNet [16] was proposed for 3D object recognition on three representations, including 3D CAD models, RGB-D data, and LiDAR point clouds. VoxNet is a kind of shallow volumetric network constructed using five layers of CDBN for real-time 3D object recognition tasks. VoxNet inputs $32 \times 32 \times 32$ voxel data. Sedaghat et al. introduced orientation-boosted voxel nets [19] to produce a correct orientation of an object and then predict the object class level as a parallel task. Recently, Wang et al. proposed NormalNet [17] with a reflection-convolution-concatenation (RCC) module. This approach can be compared to a fusion network that takes both voxels and normal vectors of voxels as inputs with resolutions of $30 \times 30 \times 30$ and fed to the network. The high-level features are generated by concatenation layer from three different poses of a 3D object. Most of the recent CNN based 3D volumetric architectures [14–16, 18, 19, 56] significantly improved the state-of-the-art performance in 3D object classification. However, to reduce the gap between 2D based CNNs and volumetric CNNs for 3D object classification, Voxception-ResNet (VRN) [56] achieved a significant improvement by utilizing inception architecture followed by Resnet [57]. The VRN is a very deep framework compared with recent volumetric networks; it contains 45 layers with 18M parameters approximately. However, voxel-based volumetric representations store both occupied and nonoccupied spaces into a data tensor and manage to feed them to a 3D CNN, which is very hard. This may limit the production of higher resolution features because of their high computational costs.

Octree-based volumetric CNNs have recently been used in various applications of 3D shape analyses because of their lossless compression landmarks with rich object boundary information in comparison to voxel grid representations [25, 58, 59]. Octree-based point-cloud compression was proposed by Schnabel et al. [60], where a point cloud is encoded based on local surface approximations that outperform previous progressive methods in terms of compression rates. To increase the geometric

resolution, Laine et al. [24] proposed a novel compression technique where the voxel data are augmented for the sparse octree to give smooth surfaces and greater geometric detail. Variational range data are fused by proposing dynamic octree partition for volume-based object reconstruction in [61]. A similar approach was proposed by Tatarchenko et al. [22]; an octree generating deep convolutional decoder was proposed to reconstruct high-resolution 3D shapes. Recently, Riegler et al. [23] proposed a novel octree representation to produce high-resolution volumetric features to feed their proposed OctNet, which was applied for 3D classification, orientation estimation, and segmentation tasks. Instead of a regular octree structure, the maximal octree depth is restricted to a depth of 3, which can be compared with $8 \times 8 \times 8$ resolutions. However, these shallow octrees are placed along regular grids. The OctNet inputs high-resolution geometric 3D volumetric data tensor consuming lower memory space, and the classification accuracy on low-scale ModelNet10 dataset improved significantly compared with other volumetric methods. However, we also considered octree representation in this paper, but we do not limit the maximal octree depth as OctNet [23] did, rather we continue the octree partition until the maximal depth is reached in accordance with regular octree representation. Our MS-VDCNN outperforms the OctNet on ModelNet [6] datasets.

3. Methods

Our proposed MS-VDCNN uses octree representation as its input. To build octree data tensors, a triangle mesh or a point cloud is preprocessed outside of the network. A list of all training models is packed with labels (a numeric number e.g. 0, 1, 2, ..., $n-1$) and their corresponding directories. The same process is also applied to preprocess the test dataset. The CNN operation on the octree volumetric data is executed using octree convolution layer, and a 2D convolution kernel with a 3D filter is applied on non-octree 3D data. Our MS-VDCNN can be performed on both CPU and GPU environments. The data preprocessing method, MS-VDCNN framework, and network operations will be described in the following.

3.1. Point-Octree Representation. Our method represents the surface geometry of a 3D object with regular point data in 3D space. For 3D CAD models (.obj/.OFF), we applied the ray shooting algorithm with the help of [62] and shot 16k rays uniformly on 3D objects. Then, we calculated the intersection of the rays on the object surface and stored the sampled dense points. In the beginning, we scaled a 3D object (a point cloud) uniformly and sent to a unit 3D bounding cube with z -axis alignment. To reduce the computational footprint, the sampled points were used to form a volumetric octree O_t structure, while the CNN operations will only be performed on nonempty octree nodes. These nonempty nodes present the boundary surface information of the 3D object. A bounding cube of a

point cloud 3D object is subdivided into eight equal pieces to construct octrees. Thus, the precious geometric properties with surface information were transformed into the smallest cube. The occupied or partially occupied nodes only undergo the recursive subdivision process. We continued this subdivision process iteratively until it reached the predefined octree depth d^{th} . The input volume resolution can be compared with octree depths, i.e., octree depths 4, 5, 6, 7, and 8, can be compared with voxel resolutions 16^3 , 32^3 , 64^3 , 128^3 , and 256^3 , respectively. The partitioning process of octrees is applied on nonempty nodes only where the resulting nodes are known as leaf or child nodes. Figure 3 depicts our octree representation of a 3D car object. The subdivision process is shown in Figure 3(a) from depth zero (d^0) to depth three (d^3). Corresponding to the original x , y , and z axis, the location of child nodes can be represented by i , j , and k axis, respectively, as follows:

$$\begin{aligned} i &= l_{n-1}2^{n-1} + \dots + l_d2^d + \dots + l_02^0; \\ j &= m_{n-1}2^{n-1} + \dots + m_d2^d + \dots + d_02^0; \\ k &= p_{n-1}2^{n-1} + \dots + p_d2^d + \dots + p_02^0; \end{aligned} \quad (1)$$

where $l, m, p \in \{0, 1\}$, d is the octree depth ($d = n-1, n-2, n-3, \dots, 1, 0$), and the values of $l_d, m_d,$ and p_d indicate the octant location in the cube at depth d can be written as

$$\xi_d = p_d2^2 + m_d2^1 + l_d2^0, \quad (2)$$

If $l_d, m_d, p_d = 1$, the point data are located in one of the odd numbering octants label, i.e., 1, 3, 5, and 7 at corresponding d^{th} depth; otherwise, the even numbering position will be discovered, i.e., 0, 2, 4, and 6 and ξ_d indicates the partition path from the origin to final octant position where the point data are located.

However, to find the sibling nodes sequentially, each newly generated leaf octants are labelled with a positive integer number κ_l . Assigning a label to nodes is a very straightforward technique, where all nonempty nodes are labelled with positive numeric numbers (e.g., 1, 2, 3, ..., l) and zero (0) for empty nodes. The total number of occupied nodes is l , and R_{\max} is the total number of leaf nodes. The labelling process is shown in Figure 3(b). The values of all labels at d^{th} depth are stored in a vector $\varphi_\xi^{d^{\text{th}}}$ can be written as

$$\varphi_\xi^{d^{\text{th}}} = \left\{ \varphi_{\xi_1}, \varphi_{\xi_2}, \varphi_{\xi_3}, \varphi_{\xi_4}, \varphi_{\xi_5}, \dots, \varphi_{(\xi(R_{\max}-1))}, \varphi_{\xi(R_{\max})} \right\}; \quad (3)$$

From equation (2) at

$$\begin{aligned} \xi_d, & \quad \text{if } l_n, m_n, p_n \neq 0, \text{ then } \varphi_\xi = \kappa_l; \\ \text{otherwise,} & \quad \varphi_\xi = 0; \end{aligned} \quad (4)$$

Finally, the averaged normal vectors are computed at the finest leaf octants as the input signal of CNN by $\gamma_n = 1/R_{\max} \sum_r^{R_{\max}} h_\kappa(r+1)$, where h_κ represents the values

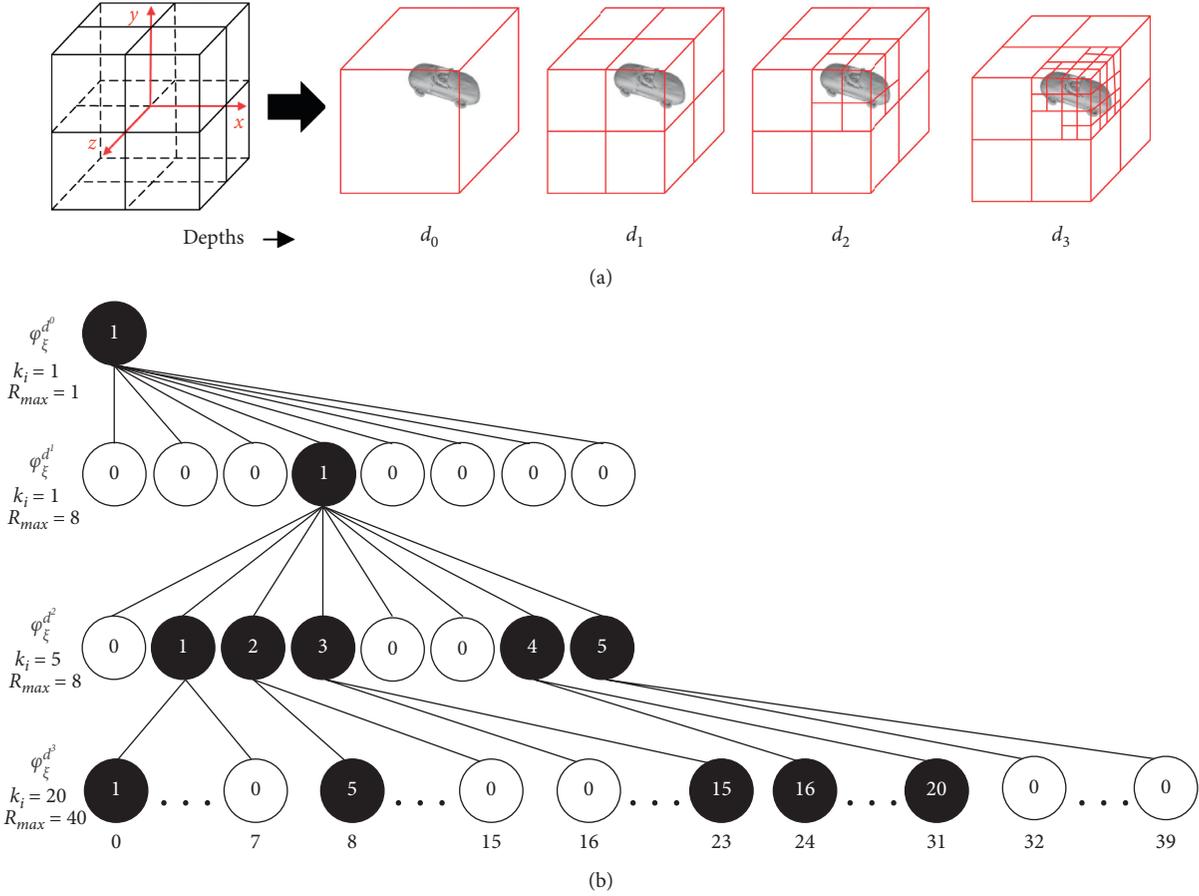


FIGURE 3: An example of our octree illustration of 3D car object: (a) octree representation according to depths and (b) labelling of nodes to build parent-child (leaf) relationship at d^{th} depth.

of leaf nodes and R_{\max} is the number of total nodes in the finest leaf. From Figure 3, we can see the octree representation of a 3D car object from depth zero (d^0) to depth three (d^3). Given a bounding box is set to d^0 , represented by a black node in Figure 3(b) at the top. Black nodes define the nodes whether it is occupied or partially occupied by the object parts. At depth d^1 , the number of occupied nodes $l = 1$, and total nodes $R_{\max} = 8$ nodes. In the following, at the second depth (d^2), the number of total nodes remain same with the number of d^1 , but R_{\max} is increased at 3rd depth ($R_{\max} = 50$) because of the number of occupied nodes increased to 5 ($l = 40$).

So, octree gives dynamic resolutions (may increase, decrease, or remain same) over the depths where memory space is increased recursively (Figure 3(b)). By contrast, voxels increase the resolution cubically and consume higher memory than the octree representation. However, it is a nontrivial task to implement a CNN operation on octree structured data because of the hierarchical tree structure that contains difference range data at d^{th} depths. We addressed this problem by keeping the same number of hidden blocks in the network with octree depths.

In addition, it is not easy to do the geometric transformation of an octree shape to perform a subvolume supervision. We need to get back the feature data into raw-

volume shapes ($x \times y \times z$). During this transformation, information loss may occur at d^{th} depth. To tackle this problem, we propose reserving the full layer of octants until a minimal number of octree depth is reached. This reservation approach may increase the size of input vectors and computational overhead slightly than regular an octree structure, but a lossless transformation would be achieved. Our method is still efficient for 3D representation at higher resolutions than full voxel methods in terms of consuming less memory space and decreasing computational footprints significantly, as described in Section 2 (Figure 2). To train multiple networks with different octree depths, a simple conditional logic function is set to define the depth number, where to reserve all octants at d^{th} depth can be written by the following equation:

$$O_{\text{TN}}^{d^{\text{th}}} = \begin{cases} 2, & \text{if } \psi_n \leq V_{R_{\max}}, \\ 3, & \text{else,} \end{cases} \quad (5)$$

where $O_{\text{TN}}^{d^{\text{th}}}$ represents the expected full-layer octree depth, ψ_n is the current octree depth in the loop, and $V_{R_{\max}}$ is the predefined input volume resolution. If we consider the volume resolution to be higher than 16^3 , then all nodes till the third depth will be stored forcibly unless the number of

the full layer values is 2 for the other $V_{R_{\max}}$. This extra reservation of nodes never affects the octree partitioning process.

To conclude, we encoded the 3D object surface with a set of points that were sorted into the nodes of octree (at the finest leaves) with a maximal number of d^{th} depths. The highest number of depths determined the precision of contour-information. Thus, the average normal vectors were computed at the finest leaf nodes, and the CNN inputs these normal vectors as the input features of the 3D object. Although the points were encoded into octree as a top-down fashion, the reverse process was applied for generating high-level features in the network. The CNN operation was optimized by using point-octree representation, and its generated high-level features will be fed further to the FC layers for the 3D object classification task. Figure 4 shows several 3D models with our point-octree representation and a comparison with corresponded original models from the ModelNet40 dataset.

3.2. Architecture of MS-VDCNN. Our MS-VDCNN is a feed-forward supervised CNN employed with multiscale residual and subvolume supervised learning tactics for the 3D object classification task.

MS-VDCNN is divided into two branches: the main branch learns global feature from the whole shape utilizing residual learning tactics and the second branch is for the subvolume supervision to learn local features from the object parts. The idea of learning from subvolume parts can be compared with a perceptual human vision system, as we also often identify the object by reviewing a particular part of it. In addition, subvolume learning works against overfitting to training data and helps to continue the learning process. Due to considering only subvolume parts of the input, it is difficult to overfit the training data. A multiscale residual block is proposed to form a high-level precise feature by combining low-level features from hierarchical levels of octree depths inside the network.

The conv. and FC are the main layers of our network. Figure 5 depicts the architecture of MS-VDCNN consisting of multi-octree-conv blocks (MOC), multiscale residual blocks (MSRB), a typical conv. layer (Tconv.), and FC layers. The MOC consists of two octree-conv layers in a series connection, and a feature channel (F_n) is increased twice to the second octree-conv. We use a common kernel size of all octree convolution layers that are set to 3 with a stride of 1. The convolution operation can be written as follows:

$$O\tau_c(\xi) = \sum_n \sum_l \sum_m \sum_p W_{\text{imp}}^n \cdot F_n^{d^{\text{th}}}(\xi_{\text{imp}}^{d^{\text{th}}}), \quad (6)$$

where W_{imp}^n is the convolution operation and $F_n^{d^{\text{th}}}$ is the n^{th} channel features vector at the d^{th} depth. The neighboring octant of ξ is represented by $\xi_{\text{imp}}^{d^{\text{th}}}$. $F_n(\cdot)$ is set to zero if $\xi_{\text{imp}}^{d^{\text{th}}}$ is not available to the corresponding octree depth, and the convolution operation will be converted to matrix product. The number of MOC will depend on the depth of octrees. For the different levels of octrees, the lengths of the MS-VDCNN will be changed. In Figure 5, the MS-VDCNN is



FIGURE 4: Our octree representation from corresponding original 3D CAD models from the ModelNet40 dataset.

designed for the fifth depth octree data, and it can be compared with $32 \times 32 \times 32$ voxels; but one unit of MOC should be added or omitted if the depth of octree is more or less 1. This MOC block will help to generate high-level discriminative feature without reducing the octree depth and keep the features invariant. Batch normalization (BN) and nonlinear rectified activation function (ReLU) are applied to the end of each octree-conv layer. The output of last conv. of the MOC block is fed to a max-pooling layer with a stride of 2. The max-pooling operation is applied on n^{th} channel features map independently, and the resolution of the featured map is downsampled by a factor of 2. The octree data tensor with d^{th} depth will be downsampled to $(d-1)^{\text{th}}$. The pooling operation $\chi_{i,j,k}^{mp}$ will select the maximum octant value from every 8 adjacent leaf octants, which can be written from equation (6) as

$$\chi_{i,j,k}^{mp} = \max_{l,m,p \in \{0,1\}} \text{pooling} \sum_o^{n-1} (F_{n-1}^d \xi_{2i+l, 2j+m, 2k+p}). \quad (7)$$

The MSRB consists of a max-pooling, an Elitwise concatenation, and ReLU layers. To put it simply, our MSRB is a fusion block (MSRB, $f(y) = f(x) + x$), where feature vectors are fused between two MOC or two conv. layers. The residual connection is illustrated in Figure 6. The pooling layer will resize the feature vector from the previous conv. layer. This adaptation is needed to match the octree depth with the current feature vector. Then, the multiscale information from multiresolution features vectors is fused by MSRM to generate high-level distinguishable features. The last output of MSRB has a size of $512 \times 8 \times 1$ corresponding to $512 \times 2 \times 2 \times 2$ voxels which can be converted by the O-V layer. A typical convolution layer is also applied to the main branch of MS-VDCNN, where the 2D convolution operation with a 3D filter (Figure 7) is used with a stride of 1. This typical convolution is commonly used for 2D image processing, whereas a 3D filter is adopted to perform on 3D volumetric data tensor. All convolution operations are computed efficiently on the GPU. We reshaped the feature vector from $(n * c * h * w)$ to $(n * (c * h * w))$ and fed it to the FC-1 directly. We have three FC layers in our main branch with lengths of 512, 512, and 40, respectively. The length of the last FC layer depends on the categories of objects in the dataset (i.e., 40 for the ModelNet40 and 10 for the ModelNet10 dataset). We applied a 50% drop out between FC-1 and FC-2. The subvolume supervision was applied to a $512 \times 2 \times 2 \times 2$ data tensor where 4D data

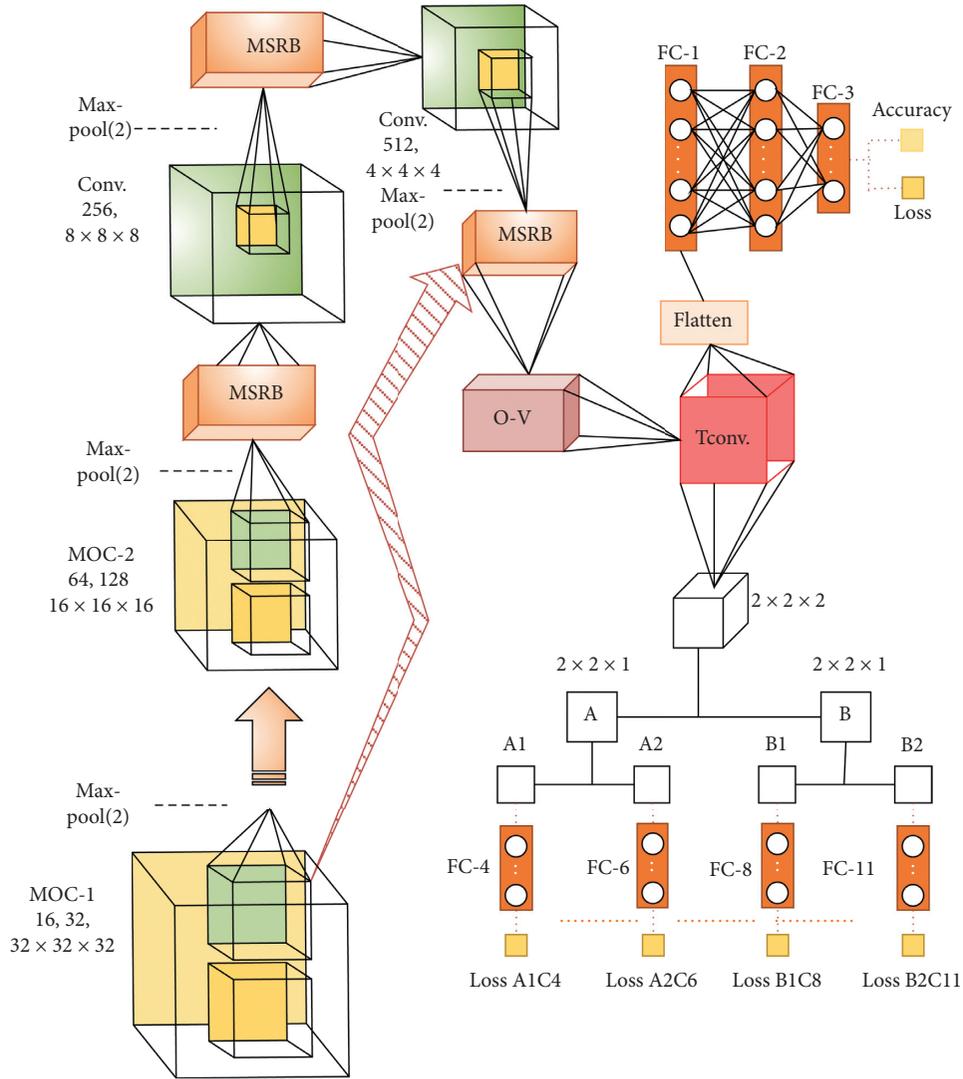


FIGURE 5: The network architecture of our proposed MS-VDCNN for an octree depth of 5.

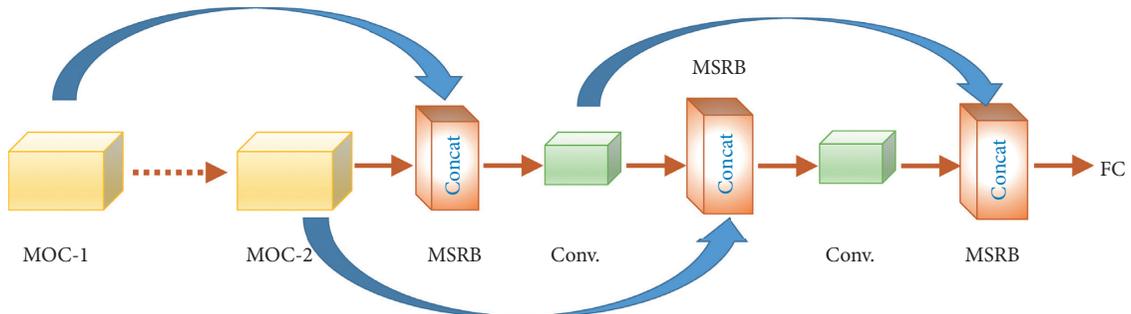


FIGURE 6: Multiscale residual block (MSRB) connection.

tensors were sliced sequentially into eight subparts. The slicing techniques are shown in Figure 5. These newly generated eight layers will be used to predict object class using a subvolume supervised learning approach. A Softmax layer was applied at the end of each FC of the sliced layers, and it simultaneously computes the average classification accuracy and loss. This training on

subvolume parts of the object can easily enjoy accuracy gains from integrated information. However, the accuracy prediction by the main branch is more authentic as it exploits whole object surface information. Finally, we collected the object class label predicted by FC-3, and the highest prediction among 40 nodes was the final class of an object.

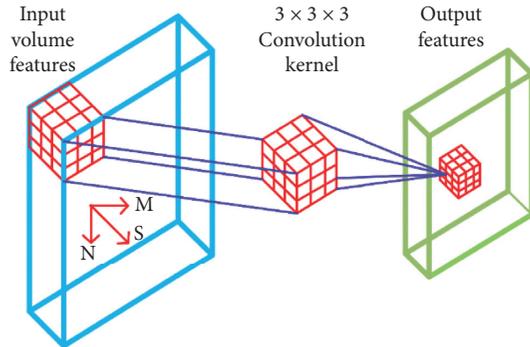


FIGURE 7: Example of 3D Convolution filter.

In addition, to evaluate the performance of our MS-VDCNN, we computed the Softmax loss quantitatively with respect to the negative log-likelihood function. The multi-class loss function can be written as

$$l(y, \hat{y}) = -\frac{1}{M} \left[\sum_1^M \sum_j^N y_j \cdot \log \left(\frac{e^{x_j}}{\sum_j^N e^{x_j}} \right) \right], \quad (8)$$

where y_j is the ground truth for the j^{th} class in N class. The $\log(\cdot)$ is the probability function of predicting the target class by the Softmax function, where x_j is the logit score of input data and \hat{x}_j is the score referred by the network for each class in N .

4. Experiments

This section demonstrates the implementation of 3D classification task by our MS-VDCNN, and it compares the results with state-of-art methods. We considered the ModelNet [6] benchmark datasets. The Caffe [63] deep learning framework was used to implement MS-VDCNN. Initially, we made 4-networks, including MS-VDCNN I, II, III, and IV, depending on the input resolutions (16^3 , 32^3 , 64^3 , and 128^3 , respectively). We trained all of the networks using an Intel Xeon-X5650 desktop PC with a Tesla K20c (5 GB) GPU.

4.1. Data Preparation. The 3D CAD models from ModelNet [6] datasets, including ModelNet40 and ModelNet10, were considered. The modelNet40 dataset contained 12311 CAD models (9843 models for the training and 2468 models for the testing) of 40 categories, and the ModelNet10 contained 4899 indoor models (908 models for the testing and 3991 for the training). The preprocessing of 3D mesh (.obj/.off) to our point-octree is described in Section 3.1. Initially, we prepared 4 sets of octree datasets from the ModelNet40 with four different octree depths 4, 5, 6, and 7, which were compared with voxel resolutions of 16^3 , 32^3 , 64^3 , and 128^3 , respectively. Before we made an octree representation, we rotate each 3D model uniformly with 12 rotations and aligned its z -axis. For a fair comparison, we follow the original training and testing split that was given in the dataset [6]. All networks are trained on the augmented training datasets ModelNet40; these had 147732 3D models

in total. The training and testing are conducted on augmented datasets that had 118116 and 29616 3D samples, respectively.

4.2. Training Scheme. We trained the 4 networks of MS-VDCNN on the abovementioned desktop PC. Our MS-VDCNN is a feed-forward supervised deep convolution neural network trained with an end-to-end fashion and optimized by the stochastic gradient descent (SGD) method with a momentum of 0.9. The initial learning rate was set to 0.1 and decreased by factor 10 after 10 epochs approximately. The batch size of all networks was set to 32 and optimized at 45 epochs approximately. Therefore, we observed the learning behavior on the training dataset and stored the weights in the interval of 5 epochs. We computed the classification accuracy and loss both on the training and the testing datasets in the interval of 2000 iterations. During the test, all augmented test samples of the prime sample were fed to MS-VDCNN in one batch, and the average classification accuracy (average instance accuracy) was computed.

At first, we trained our networks on the Model40 dataset using 12 rotations of each 3D object around the horizontal axis and observed the influence of the octree depths. After we got the best octree depth, we prepared several datasets by augmenting each 3D object with a set of rotations and recorded the effective classification accuracy. We applied the transfer learning approach on ModelNet10 and used a pretrained model from ModelNet40.

4.2.1. 3D Object Classification and Result Analysis. To perform 3D object classification, we trained our networks from scratch within the ModelNet40 dataset. Similar to volumetric networks 3DShapeNets [15] and VoxNet [16], we augmented the dataset with 12 poses of each 3D object and varied the input resolutions from 16^3 to 128^3 . We focused on improving the classification accuracy and observed the impact at higher resolutions. MS-VDCNN-I contains one MOC block, and the rest of the blocks are similar to the network is shown in Figure 5. The number of MOC blocks was increased by one block by increasing twice the input resolution. An increase in octree depths is the technique of increasing the volume resolutions for octree representation. So, MS-VDCNN-I had one MOC block for depth 4, MS-VDCNN-II had two MOC blocks for depth 5, MS-VDCNN-III has four MOC blocks for depth 6, and MS-VDCNN-IV had five MOC blocks for depth 7. Each MOC block contains two octree-conv. and one max-pooling layer.

The training and the testing accuracy over the epochs is demonstrated in Figure 8. During the training, a single epoch contained 118116 samples that are needed to be trained. We set our batch size to 32 samples for all networks, and our small GPU with a 5 GB memory could easily compute them. In Figures 8(a)–8(d), the classification accuracy map during the training was drawn by our four MS-VDCNNs, and the comparison of test accuracy among them is plotted in Figure 9(a). We saw the highest classification accuracy was achieved by MS-VDCNN-IV (0.9233), which was designed with octree depth 7; it can be compared to a 128^3 voxel resolution. In Figure 8, all the networks are

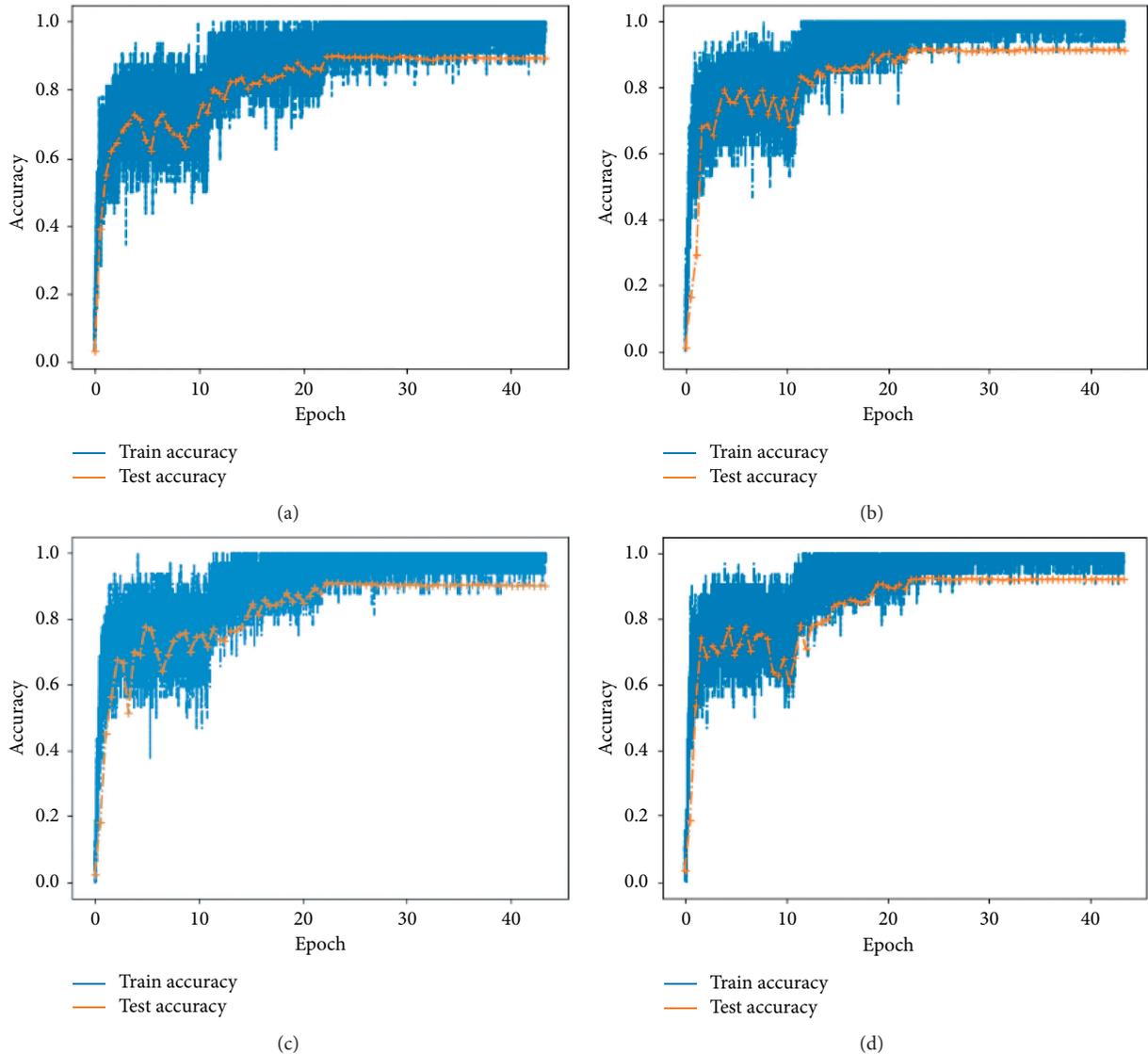


FIGURE 8: Classification accuracy during the training on 12 views: (a) MS-VDCNN-I (D-4), (b) MS-VDCNN-II (D-5), (c) MS-VDCNN-III (D-6), and (d) MS-VDCNN-IV (D-7).

almost optimized just after 20 epochs, and we ended the training at 40 epochs, approximately.

In our volumetric CNN, the volume resolution is downsampled after each CNN operation by a max-pooling layer where our MSRB blocks produced high-level features by concatenating features from previous CNN outputs. This process will be continued through all the convolution blocks until it reaches the FC. Our subvolume learning scheme helped to increase the predicting probability to the actual observation label by learning information from subparts of an object. This tactic might help to generate a more accurate distinguishable prediction score of the desired class.

Figure 9(a) depicts the test’s accuracy achieved by our MS-VDCNNs over the epochs and a comparison with other two state-of-the-art volumetric methods (3DShapeNets and VoxNet) for 3D object classification. From the accuracy plot, it can be seen that our MS-VDCNNs outperform under all

resolutions, even with lower input resolutions (16^3), compared with the 32^3 used by VoxNet and 3DShapeNets. The best classification accuracy of 92.33 % was achieved by our MS-VDCNN-IV (D-7) while using input resolution of 128^3 . Figure 9(b) shows the learning behavior during the training by our MS-VDCNN over the epochs, and it nearly reaches zero (0) after 40 epochs where we decided to stop the training.

However, at depth 6, the increasing rate of classification accuracy is comparatively slowed. This minor ambiguity can be found between some classes of objects at higher resolutions because of visual similarities [23]. This issue can be successfully resolved by using multiscale learning [64], as it combines both low- and high-level features. In Table 1, we can see that the sample accuracy is increased by 0.3% at depth 5 with multiscale learning, whereas the accuracy is greatly improved at depth 6 (Figure 10). In addition, our

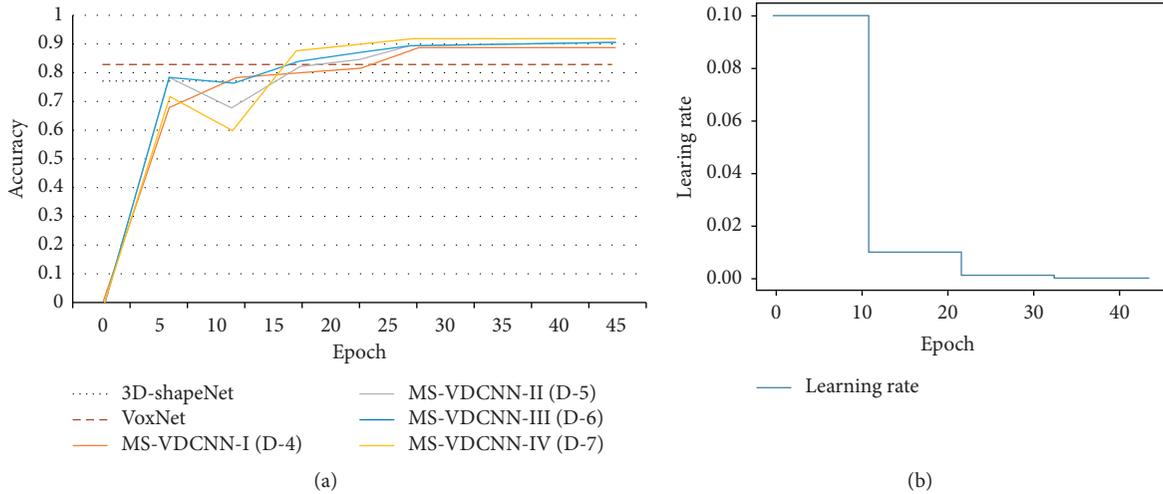


FIGURE 9: A comparison of classification accuracy on ModelNet 40: (a) our CNNs’ performance over the epochs, MS-VDCNN-I, MS-VDCNN-II, MS-VDCNN-III, and MS-VDCNN-IV is designed for input resolutions 16^3 , 32^3 , 64^3 , and 128^3 , respectively; and a maximum classification accuracy of VoxNet and 3DShapeNets is shown with a static value collected from the corresponding papers. Both VoxNet and 3DShapeNets used 32^3 voxel grids and (b) a learning rate map of our MS-VDCNNs over the epochs.

TABLE 1: A comparison between our octree with voxel representation with 32^3 volume resolution under the same network (MS-VDCNN-II) on the Modelnet40 dataset.

Representation method		Multiscale	Accuracy (%)
Volume resolution 32^3	Octree	No	91.1
	Voxel	No	90.3
	Octree	Yes	91.4
	Voxel	Yes	90.5

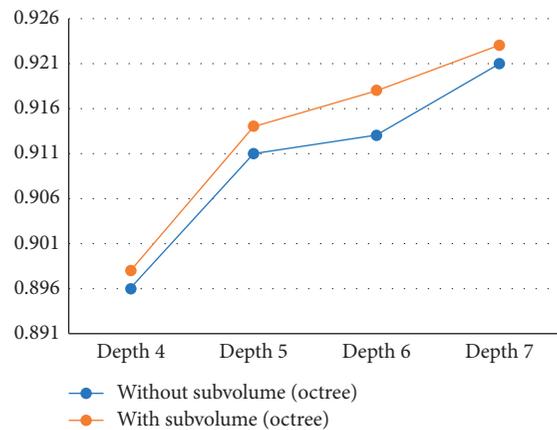


FIGURE 10: Classification accuracy of various octree depths with the effect of subvolume learning on ModelNet40 dataset.

proposed method improved the sample accuracy by 0.9% than the voxel representation. There is some redundant information in voxel representation, as it encodes both occupied and nonoccupied spaces, but the octree encodes very precise information by encoding only occupied spaces in the volume. In addition, to evaluate the performance of our framework, the log-loss was computed by using equation (8) to measure the uncertainty of the prediction. If the predicted probability increases, then the log-loss function

should be decreased. Theoretically, a log-loss of zero (0) would be an indicator for a perfect CNN framework.

It also indicates how the parameters are chosen to generalize the network. Figure 11 shows the logarithmic loss was computed by our MS-VDCNNs (on ModelNet40 dataset) which were used to generate the accuracy plots (Figure 8). We can compare the relationship between accuracy and loss plots from Figures 8 and 11. Log-loss decreases (Figure 9) as the predicted probability converges

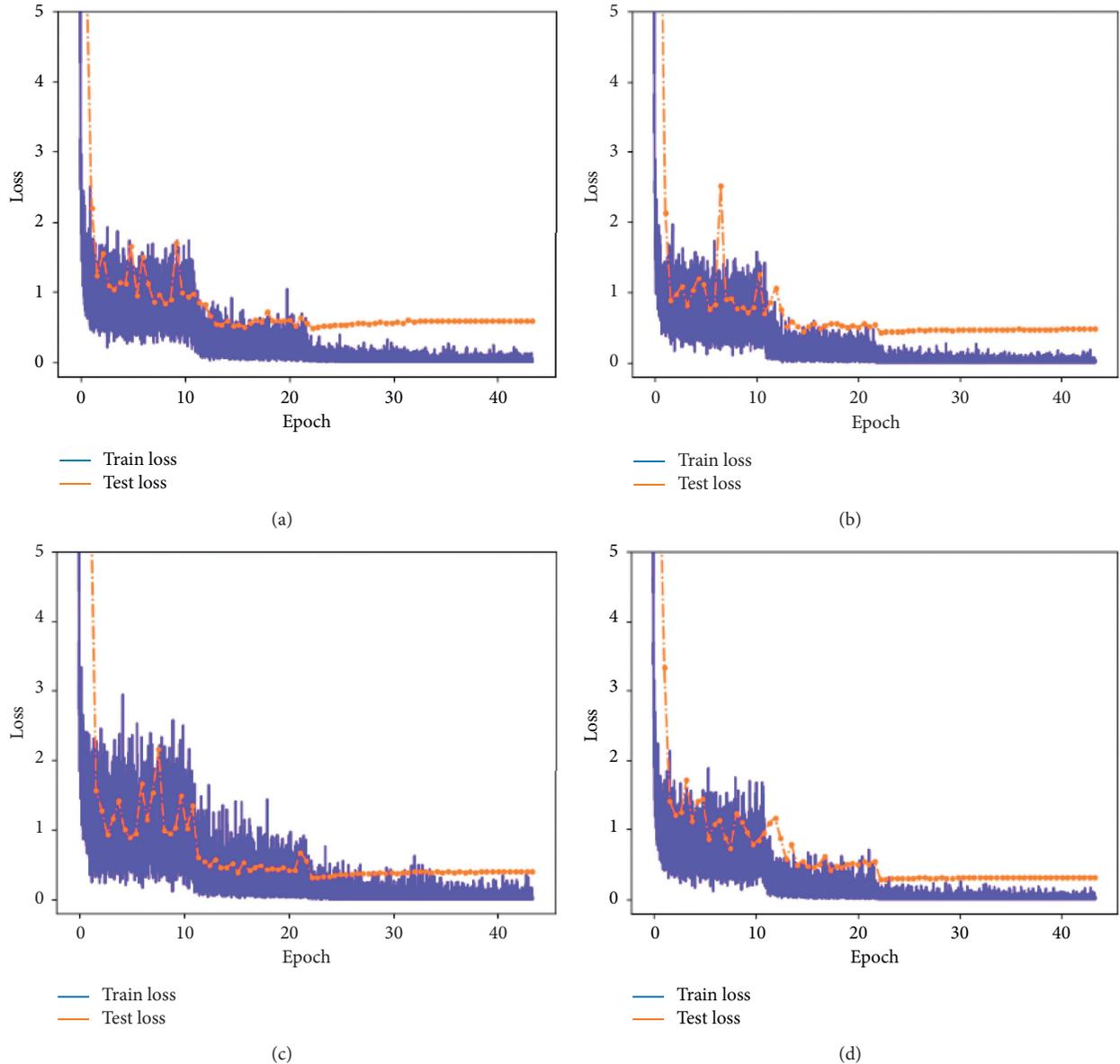


FIGURE 11: Log-loss plots during the training on 12 views: (a) MS-VDCNN-I, (b) MS-VDCNN-II (c) MS-VDCNN-III, and (d) MS-VDCNN-IV.

from the actual label resulting accuracy gains over the epochs shown in Figure 8.

However, we conducted several experiments on several augmented datasets to observe the effect on samples' numbers in different views. In general, the viewpoints of a 3D object strongly depend on their orientations where the surface and boundary information are powerful features for classification application. We applied several rotations of each object in its azimuth axis only, and we prepared 6-augmented dataset from ModelNet40 view 3, 6, 9, 12, 24, and 30, respectively. We trained MS-VDCNN-IV independently on all of the datasets. In general, DL requires more samples to be trained, and increasing the views is the most common practice to improve the accuracy [17, 20, 48]. However, too many views may produce wrong category information

among similar appearance objects. Figure 12 demonstrates the accuracy plots over the epochs during the training on augmented datasets. Our networks also optimized noticeably faster on large scale datasets, and a small number of epochs were required to be trained (Figures 12(a)–12(f)). The best accuracy so far achieved by our MS-VDCNN is 92.93% using 24 views (Figure 13), where 0.6% accuracy increased in comparison to 12-view, and log-loss was computed to 0.32 in Figure 14. However, 0.15% accuracy loss was obtained for 30-view than 24-view, and we assumed this may happened because of the insufficient number of invariant samples per class.

Table 2 shows a comparison result between our methods and other recent approaches including volumetric, multi-view, and point representation inputs. Overall, the 3D object

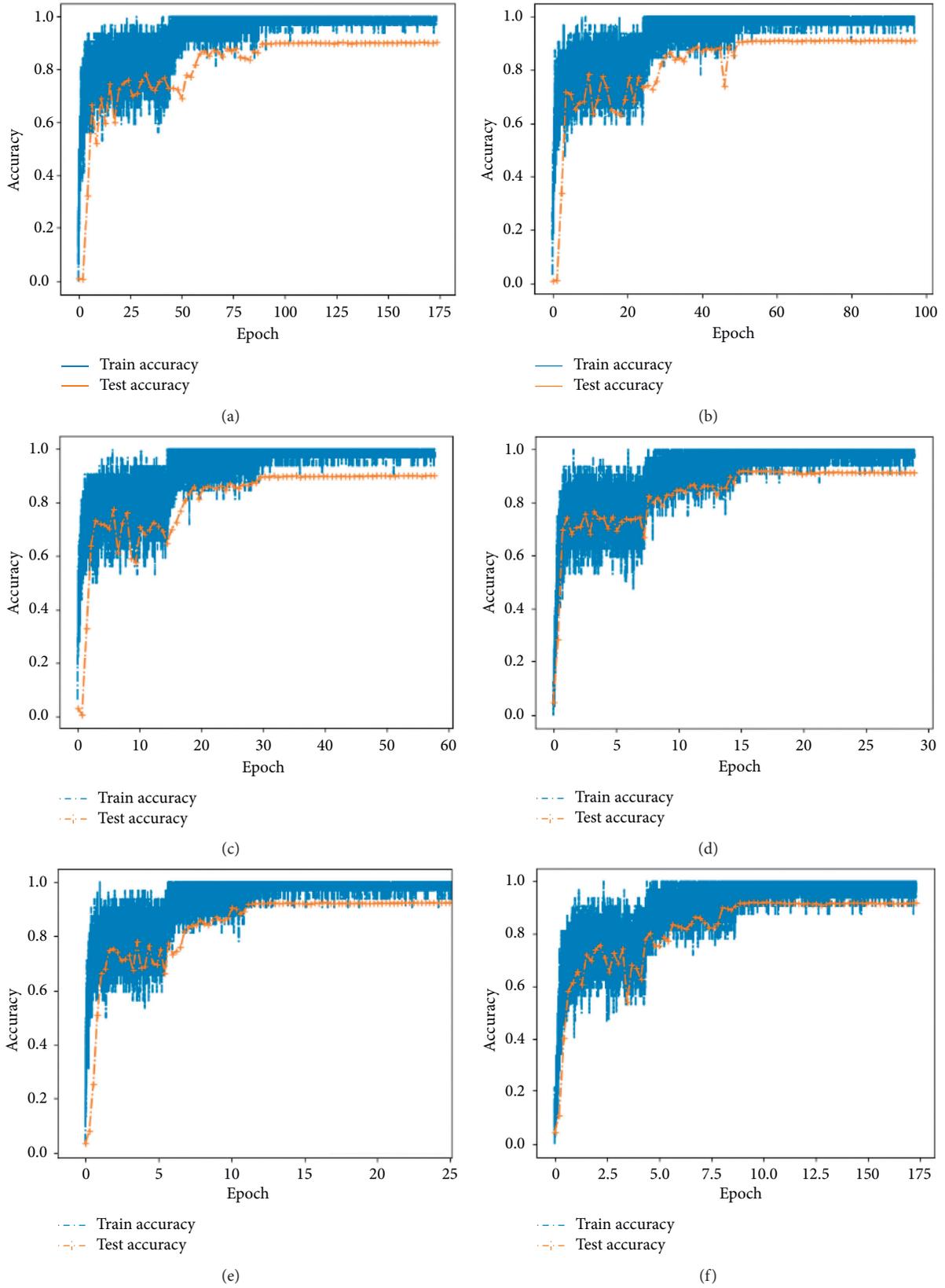


FIGURE 12: Classification accuracy plots by MS-VDCNN-IV (128³): (a) view-3, (b) view-6, (c) view-9, (d) view-12, (e) view-24, and (f) view-30.

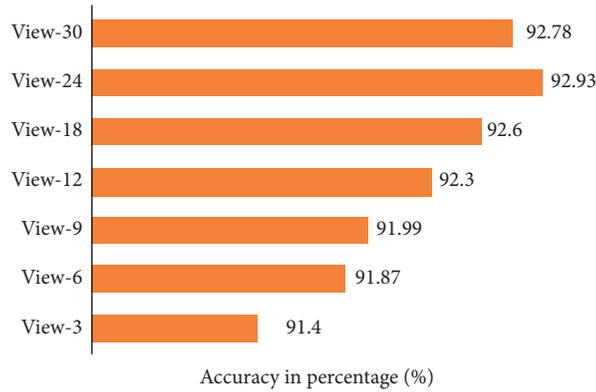


FIGURE 13: Comparison of classification accuracy on different augmented datasets of ModelNet40.

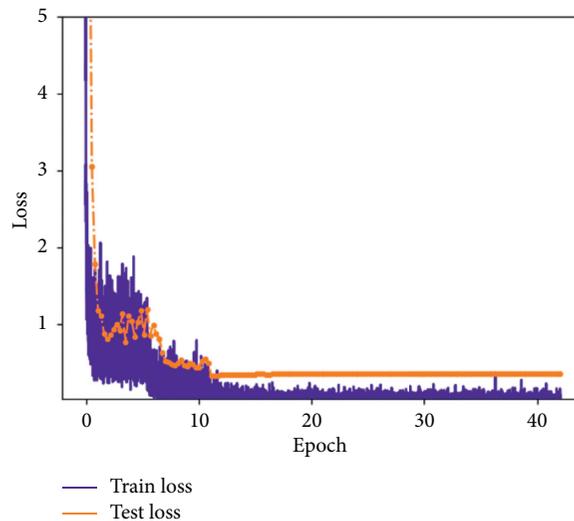


FIGURE 14: Epoch vs loss is computed on 24 views (0.32).

classification performance by multiview and point cloud CNNs is better than volumetric CNNs, except VRN ensemble [56]. Our proposed network outperforms among single types volumetric approaches on the ModelNet40 dataset with a small number of parameters (6.2 M approximately) shown in Table 2. The total parameter of NormalNet (6.5 M) [17] is close to ours, but the performance of this network is worse than our model (91.9 % vs. 92.93%). We also tested our MS-VDCNN on a low scale ModelNet10 dataset. It achieved an accuracy of 95.3% when using a pretrained model from the ModelNet40, which is also best in other volumetric 3D CNNs, reported in Table 2, except for the VRN ensembles [56]. Our MS-VDCNN achieved slightly better performance than the single VRN model (about 18 M parameters) [56] on both the ModelNet40 and ModelNet10; therefore, the performance of VRN ensembles may not apply in general cases [56].

To conclude, our proposed MS-VDCNN method is highly capable to run on a small GPU capacity (5 GB) at higher resolutions which could be an option for real-time operation, because the higher resolution might give more

discriminative feature among objects with similar appearances [23]. In addition, MS-VDCNN can also perform significantly better on a small-scale dataset (i.e., ModelNet10) where an increase in the number of views and resolutions improve the performance.

To measure the effectiveness of our proposed MS-VDCNN on point-octree, we compared the memory consumption and the average time for one complete pass. We ran 500 forward and backward passes on the GPU and computed how much time and memory were consumed (plotted in Figure 2). The batch size was 32. Our MS-VDCNN with point-octree worked faster and consumed less memory under all resolutions when the input resolution was equal or higher than 16^3 ($\geq 16^3$) in comparison to the full voxel method. Due to our GPU memory limitation (5 GB), we cannot proceed on the higher resolutions above 32^3 by the full voxel method, where we show the anticipated estimation using dotted lines in the plots (Figures 2(a) and 2(b) for memory and runtime, respectively). However, our GPU easily fits up to 128^3 input resolutions using point-octree in an MS-VDCNN. These

TABLE 2: Classification results by different methods on ModelNet40 and ModelNet10 datasets.

Network type	Method	Input	Size	Pretrain	Augmentation	ModelNet40 (%)	ModelNet10 (%)
Single volumetric	MS-VDCNN(Ours)	Volumetric	6.2 M	ModelNet40	24	92.93	95.3
	3DShapeNets [15]	Volumetric	38 M	ModelNet40	12	77.32	83.5
	VoxNet [16]	Volumetric	0.92 M	—	12	83	92
	Voxception [56]	Volumetric	—	ModelNet40	24	90	93.28
	OctNet [23]	Volumetric	—	—	—	86.5	90.9
	VRN [56]	Volumetric	18 M	ModelNet40	24	91.33	93.61
	Aniprobng [14]	Volumetric	—	—	60	85.6	—
Ensemble volumetric	NormalNet [17]	Volumetric + norm: vector	6.5 M	—	20	88.8	93.1
	VRN ensemble [56]	Volumetric	90 M	ModelNet40	24	95.54	97.14
Point cloud	FusionNet [20]	Volumetric + multiview	118 M	ImageNet ModelNet40	60	90.80	93.1
	PointNet [50]	Points	0.45 M	—	—	86.2	89.2
	PointNet++ [52]	Points	—	—	—	—	90.7
	3D Capsule [55]	Points	—	—	—	92.7	94.7
	RS-CNN [54]	Points	1.41	—	—	93.6	—
Multiview	MVCNN [11]	Multiview	—	ImageNet	80	90.10	—
	Ma et al. [46]	Multiview	—	ImageNet	12	91.05	95.29
	3D2SeqViews [47]	Multiview	—	ImageNet	12	93.40	94.71

results demonstrate the effectiveness of our proposed method that runs approximately seven times faster and consumed four times less memory than the full voxel method at the higher resolutions.

5. Conclusions and Future Works

We presented a novel point-octree-based volumetric deep convolutional neural network for 3D object classification. We took into account the challenges of 3D object representation to CNN and introduced an optimal point-octree representation to our proposed MS-VDCNN where multi-scale hierarchical and subvolume learning tactics improved the performance of the network. Our proposed 3D deep learning framework focused on both global and local features, and it learned sparse geometric structures from both a full and partial part of the 3D object. The experimental results revealed that increases in input resolutions and the training samples have some effect to boost up the classification accuracy. Our method significantly improved the classification performance on ModelNet datasets, and it runs seven times faster and consumed four times less memory, approximately, in comparison to full voxel methods. In future work, we would like to extend our experiments to other 3D databases and investigate on other 3D computer vision problems where learning features at higher resolutions is crucially important such as multiview 3D reconstructions and 3D appearance analyses.

Data Availability

A public open-source ModelNet datasets were used to support this study and are available at arXiv.org; csgt;

[arXiv:1512.03012](http://arXiv.org). The source of datasets is cited at relevant places within the text as references in Section 4.1 and cited in [6].

Conflicts of Interest

The authors declare no conflicts of interest.

Acknowledgments

This study was partially supported by the project of the Shanghai Science and Technology Committee (no. 18510760300), Anhui Natural Science Foundation (grant no. 1908085MF178), and Anhui Excellent Young Talents Support Program Project (grant no. gxyqZD2019069).

References

- [1] E. Ahmed, A. Saint, A. E. R. Shabayek, K. Cherenkova, and D. Aouada, “Deep learning advances on different 3D data representations: a survey,” 2019, <http://arxiv.org/abs/1808.01462v2>.
- [2] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “ImageNet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., pp. 1097–1105, Curran Associates, Inc., Red Hook, NY, USA, 2012.
- [3] C. Szegedy, W. Liu, Y. Jia et al., “Going deeper with convolutions,” February 2019, <http://arxiv.org/abs/1409.4842>.
- [4] W. Wei, T. Can, W. Xin, L. Yanhong, H. Yongle, and L. Ji, “Image object recognition via deep feature-based adaptive joint sparse representation,” *Computational Intelligence and Neuroscience*, vol. 2019, Article ID 8258275, 9 pages, 2019.
- [5] A. Voulodimos, N. Doulamis, A. Doulamis, and E. Protopapadakis, “Deep learning for computer vision: a brief

- review,” *Computational Intelligence and Neuroscience*, vol. 2018, Article ID 7068349, 13 pages, 2018.
- [6] A. X. Chang, T. Funkhouser, L. Guibas et al., “ShapeNet: an information-rich 3D model repository,” February 2019, <http://arxiv.org/abs/1512.03012>.
 - [7] S. Choi, Q.-Y. Zhou, S. Miller, and V. Koltun, “A large dataset of object scans,” May 2019, <http://arxiv.org/abs/1602.02481>.
 - [8] S. Song, S. P. Lichtenberg, and J. Xiao, “SUN RGB-D: A RGB-D scene understanding benchmark suite,” in *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 567–576, Boston, MA, USA, June 2015.
 - [9] I. K. Kazmi, L. You, and J. J. Zhang, “A survey of 2D and 3D shape descriptors,” in *Proceedings of the 2013 10th International Conference Computer Graphics, Imaging and Visualization*, pp. 1–10, Los Alamitos, CA, USA, August 2013.
 - [10] R. Socher, B. Huval, B. Bath, C. D. Manning, and A. Y. Ng, “Convolutional-recursive deep learning for 3D object classification,” in *Proceedings of the 25th International Conference on Neural Information Processing Systems*, pp. 656–664, Lake Tahoe, NV, USA, December 2012.
 - [11] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, “Multi-view convolutional neural networks for 3D shape recognition,” December 2018, <https://arxiv.org/abs/1505.00880>.
 - [12] Z. Xie, K. Xu, W. Shan, L. Liu, Y. Xiong, and H. Huang, “Projective feature learning for 3D shapes with multi-view depth images,” *Computer Graphics Forum*, vol. 34, no. 7, pp. 1–11, 2015.
 - [13] B. Shi, S. Bai, Z. Zhou, and X. Bai, “DeepPano: deep panoramic representation for 3-D shape recognition,” *IEEE Signal Processing Letters*, vol. 22, no. 12, pp. 2339–2343, 2015.
 - [14] C. R. Qi, H. Su, M. NieBner, A. Dai, M. Yan, and L. J. Guibas, “Volumetric and multi-view CNNs for object classification on 3D data,” in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5648–5656, Las Vegas, NV, USA, June 2016.
 - [15] Z. Wu, Shuran Song, Aditya Khosla et al., “D ShapeNets: a deep representation for volumetric shapes,” 3, in *Proceedings of the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 1912–1920, Boston, MA, USA, June 2015.
 - [16] D. Maturana and S. Scherer, “VoxNet: a 3D convolutional neural network for real-time object recognition,” in *Proceedings of the 2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 922–928, Hamburg, Germany, September 2015.
 - [17] C. Wang, M. Cheng, F. Sohel, M. Bennamoun, and J. Li, “NormalNet: a voxel-based CNN for 3D object classification and retrieval,” *Neurocomputing*, vol. 323, pp. 139–147, 2019.
 - [18] S. Zhi, Y. Liu, X. Li, and Y. Guo, “Toward real-time 3D object recognition: a lightweight volumetric CNN framework using multitask learning,” *Computers & Graphics*, vol. 71, pp. 199–207, 2018.
 - [19] N. Sedaghat, M. Zolfaghari, E. Amiri, and T. Brox, “Orientation-boosted voxel nets for 3D object recognition,” November 2018, <http://arxiv.org/abs/1604.03351>.
 - [20] V. Hegde and R. Zadeh, “FusionNet: 3D Object Classification Using Multiple Data Representations,” March 2019, <http://arxiv.org/abs/1607.05695>.
 - [21] W. Kehl, T. Holl, F. Tombari, S. Ilic, and N. Navab, “An octree-based approach towards efficient variational range data fusion,” August 2016, <http://arxiv.org/abs/1608.07411>.
 - [22] M. Tatarchenko, A. Dosovitskiy, and T. Brox, “Octree generating networks: efficient convolutional architectures for high-resolution 3D outputs,” February 2019, <http://arxiv.org/abs/1703.09438>.
 - [23] G. Riegler, A. O. Ulusoy, and A. Geiger, “OctNet: Learning Deep 3D Representations at High Resolutions,” in *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, Honolulu, HI, USA, July 2017.
 - [24] S. Laine and T. Karras, “Efficient sparse voxel octrees,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 8, pp. 1048–1059, 2011.
 - [25] A. Miller, V. Jain, and J. L. Mundy, “Real-time rendering and dynamic updating of 3-d volumetric data,” in *Proceedings of the Fourth Workshop on General Purpose Processing on Graphics Processing Units-GPGPU-4*, p. 1, Newport Beach, CA, USA, March 2011.
 - [26] F. Calakli and G. Taubin, “SSD: smooth signed distance surface reconstruction,” *Computer Graphics Forum*, vol. 30, no. 7, pp. 1993–2002, 2011.
 - [27] S. Fuhrmann and M. Goesele, “Fusion of depth maps with multiple scales,” in *Proceedings of the 2011 SIGGRAPH Asia Conference*, pp. 148:1–148:8, Hong Kong, Hong Kong, December 2011.
 - [28] B. Ummenhofer and T. Brox, “Global, dense multiscale reconstruction for a billion points,” *International Journal of Computer Vision*, vol. 125, no. 1–3, pp. 82–94, 2017.
 - [29] F. Steinbrücker, J. Sturm, and D. Cremers, “Volumetric 3D mapping in real-time on a CPU,” in *Proceedings of the 2014 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2021–2028, Hong Kong, Hong Kong, May 2014.
 - [30] S. Har-Peled, ‘Quadrees—Hierarchical Grids’ *Geometric Approximation Algorithms*, Vol. 173, American Mathematical Society, Providence, RI, USA, 2011.
 - [31] J. Ni, T. Gong, Y. Gu, J. Zhu, and X. Fan, “An improved deep residual network-based semantic simultaneous localization and mapping method for monocular vision robot,” *Computational Intelligence and Neuroscience*, <https://www.hindawi.com/journals/cin/2020/7490840/>, May 2020.
 - [32] Y. Guo, M. Bennamoun, F. Sohel, M. Lu, J. Wan, and N. M. Kwok, “A comprehensive performance evaluation of 3D local feature descriptors,” *International Journal of Computer Vision*, vol. 116, no. 1, pp. 66–89, 2016.
 - [33] R. Rostami, F. S. Bashiri, B. Rostami, and Z. Yu, “A survey on data-driven 3D shape descriptors,” *Computer Graphics Forum*, vol. 38, no. 1, pp. 356–393, 2019.
 - [34] Z. Liu, S. Chen, S. Bu, and K. Li, “High-level semantic feature for 3D shape based on deep belief networks,” in *Proceedings of the 2014 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 1–6, Chengdu, China, July 2014.
 - [35] Z. Han, Z. Liu, C.-M. Vong et al., “Deep spatiality: unsupervised learning of spatially-enhanced global and local 3D features by deep neural network with coupled Softmax,” *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 3049–3063, 2018.
 - [36] W. Huang, B. Lai, W. Xu, and Z. Tu, “3D volumetric modeling with introspective neural networks,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, pp. 8481–8488, 2019.
 - [37] J. Xie, Z. Zheng, R. Gao, W. Wang, S.-C. Zhu, and Y. N. Wu, “Learning Descriptor Networks for 3D Shape Synthesis and Analysis,” in *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8629–8638, Salt Lake City, UT, USA, June 2018.
 - [38] Z. Han, Z. Liu, J. Han, C. Vong, S. Bu, and C. L. P. Chen, “Unsupervised learning of 3-D local features from raw voxels based on a novel permutation voxelization strategy,” *IEEE Transactions on Cybernetics*, vol. 49, no. 2, pp. 481–494, 2019.
 - [39] Z. Han, Z. Liu, J. Han, C.-M. Vong, S. Bu, and C. L. P. Chen, “Mesh convolutional restricted Boltzmann machines for

- unsupervised learning of features with structure preservation on 3-D meshes,” *IEEE Transactions on Neural Networks and Learning Systems*, vol. 28, no. 10, pp. 2268–2281, 2017.
- [40] Z. Han, Z. Liu, J. Han, C.-M. Vogt, S. Bu, and X. Li, “Unsupervised 3D local feature learning by circle convolutional restricted Boltzmann machine,” *IEEE Transactions on Image Processing*, vol. 25, no. 11, pp. 5331–5344, 2016.
- [41] P. Papadakis, I. Pratikakis, T. Theoharis, and S. Perantonis, “PANORAMA: a 3D shape descriptor based on panoramic views for unsupervised 3D object retrieval,” *International Journal of Computer Vision*, vol. 89, no. 2-3, pp. 177–192, 2010.
- [42] Z. Cao, Q. Huang, and R. Karthik, “3D object classification via spherical projections,” in *Proceedings of the 2017 International Conference on 3D vision (3DV)*, pp. 566–574, Qingdao, China, October 2017.
- [43] A. Sinha, J. Bai, and K. Ramani, “Deep learning 3D shape surfaces using geometry images,” in *Computer Vision—ECCV 2016*, B. Leibe, J. Matas, N. Sebe, and M. Welling, Eds., vol. 9910, pp. 223–240, Springer International Publishing, Cham, Switzerland, 2016.
- [44] O. Russakovsky, J. Deng, H. Su et al., “ImageNet large scale visual recognition challenge,” *International Journal of Computer Vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [45] E. Johns, S. Leutenegger, and A. J. Davison, “Pairwise decomposition of image sequences for active multi-view recognition,” in *Proceedings of the 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3813–3822, Las Vegas, NV, USA, June 2016.
- [46] C. Ma, Y. Guo, J. Yang, and W. An, “Learning multi-view representation with LSTM for 3-D shape recognition and retrieval,” *IEEE Transactions on Multimedia*, vol. 21, no. 5, pp. 1169–1182, May 2019.
- [47] Z. Han, H. Lu, Z. Liu et al., “3D2SeqViews: aggregating sequential views for 3D global feature learning by CNN with hierarchical attention aggregation,” *IEEE Transactions on Image Processing*, vol. 28, no. 8, pp. 3986–3999, 2019.
- [48] J. Zhao, X. Xie, X. Xu, and S. Sun, “Multi-view learning overview: recent progress and new challenges,” *Information Fusion*, vol. 38, pp. 43–54, 2017.
- [49] R. Klokov and V. Lempitsky, “Escape from cells: deep Kd-networks for the recognition of 3D point cloud models,” in *Proceedings of the 2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 863–872, Venice, Italy, October 2017.
- [50] R. Q. Charles, H. Su, M. Kaichun, and L. J. Guibas, “PointNet: Deep learning on point sets for 3D classification and segmentation,” in *Proceedings of the 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 77–85, Honolulu, HI, USA, July 2017.
- [51] Y. Yang, C. Feng, Y. Shen, and D. Tian, “FoldingNet: point cloud auto-encoder via deep grid deformation,” in *Proceedings of the 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 206–215, June 2018.
- [52] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, “PointNet++: deep hierarchical feature learning on point sets in a metric space,” February 2019, <http://arxiv.org/abs/1706.02413>.
- [53] H. You, Y. Feng, R. Ji, and Y. Gao, “PVNet: a joint convolutional network of point cloud and multi-view for 3D shape recognition,” 2019, <http://arxiv.org/abs/1808.07659>.
- [54] Y. Liu, B. Fan, S. Xiang, and C. Pan, “Relation-shape convolutional neural network for point cloud analysis,” in *Proceedings of the 2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Long Beach, CA, USA, June 2019.
- [55] A. Cheraghian and L. Petersson, “3DCapsule: extending the Capsule architecture to classify 3D point clouds,” in *Proceedings of the 2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pp. 1194–1202, Village, HI, USA, January 2019.
- [56] A. Brock, T. Lim, J. M. Ritchie, and N. Weston, “Generative and Discriminative Voxel Modeling with Convolutional Neural Networks,” February 2019, <http://arxiv.org/abs/1608.04236>.
- [57] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” February 2019, <http://arxiv.org/abs/1512.03385>.
- [58] D. Meagher, “Geometric modeling using octree encoding,” *Computer Graphics and Image Processing*, vol. 19, no. 2, pp. 129–147, 1982.
- [59] D. Meagher, “Octree encoding: a new technique for the representation, manipulation and display of arbitrary 3-d objects by computer,” Technical Report IPL-TR-80-111, Rensselaer Polytechnic Institute, Troy, NY, USA, 1980.
- [60] R. Schnabel and R. Klein, “Octree-based point-cloud compression,” in *Proceedings of the Eurographics Symposium on Point-Based Graphic*, Boston, MA, USA, 2006.
- [61] W. Kehl, T. Holl, F. Tombari, S. Ilic, and N. Navab, “An octree-based approach towards efficient variational range data fusion,” in *Proceedings of the British Machine Vision Conference 2016*, pp. 21.1–21.12, York, UK, September 2016.
- [62] P. Alliez, S. Tayeb, and C. Wormser, *CGAL 4.14—3D Fast Intersection and Distance Computation (AABB Tree): User Manual*, 4.14, 2019, <https://doc.cgal.org/4.14.3/Manual/packages.html#PkgAABBTree>.
- [63] Y. Jia, E. Shelhamer, J. Donahue et al., “Caffe: convolutional architecture for fast feature embedding,” in *Proceedings of the ACM multimedia (ACMMM)*, pp. 675–678, Orlando, FL, USA, November 2014.
- [64] Y. Du, Y. Fan, X. Liu, Y. Luo, J. Tang, and P. Liu, “Multiscale cooperative differential evolution algorithm,” *Computational Intelligence and Neuroscience*, vol. 2019, Article ID 5259129, 17 pages, 2019.