

VERAM: View-Enhanced Recurrent Attention Model for 3D Shape Classification

Songle Chen, Lintao Zheng, Yan Zhang, Zhixin Sun and Kai Xu*

Abstract—Multi-view deep neural network is perhaps the most successful approach in 3D shape classification. However, the fusion of multi-view features based on max or average pooling lacks a view selection mechanism, limiting its application in, e.g., multi-view active object recognition by a robot. This paper presents VERAM, a view-enhanced recurrent attention model capable of actively selecting a sequence of views for highly accurate 3D shape classification. VERAM addresses an important issue commonly found in existing attention-based models, i.e., the unbalanced training of the subnetworks corresponding to next view estimation and shape classification. The classification subnetwork is easily overfitted while the view estimation one is usually poorly trained, leading to a suboptimal classification performance. This is surmounted by three essential *view-enhancement strategies*: 1) enhancing the information flow of gradient backpropagation for the view estimation subnetwork, 2) devising a highly informative reward function for the reinforcement training of view estimation and 3) formulating a novel loss function that explicitly circumvents view duplication. Taking grayscale image as input and AlexNet as CNN architecture, VERAM with 9 views achieves instance-level and class-level accuracy of 95.5% and 95.3% on ModelNet10, 93.7% and 92.1% on ModelNet40, both are the state-of-the-art performance under the same number of views.

Index Terms—3D shape classification, multi-view 3D shape recognition, visual attention model, recurrent neural network, reinforcement learning, convolutional neural network.



1 INTRODUCTION

3D shape classification is a fundamental problem in the field of computer graphics and computer vision. It finds applications from traditional computer aided design and medical imaging to cutting-edge mixed reality and robot navigation. The challenge of 3D shape classification stems from the difficulty of characterizing 3D surface geometry, the variety of 3D transformation and deformation, and the imperfection of geometry and/or topology, etc. Hundreds of hand-crafted 3D shape descriptors have been proposed, either from 2D rendered views [1], [2], [3], [4], [5] or directly on 3D models [6], [7], [8], [9]. They are, however, often carefully designed to characterize only one or a few aspects of 3D shapes, making them hard to generalize well.

Recently, inspired by the advances in image classification using convolutional neural networks (CNN) [10], [11], multi-view CNN (MVCNN) was presented for 3D shape classification [12], [13], [14]. By leveraging massive image databases such as ImageNet [15] to pre-train the CNN and learn image descriptors for general vision tasks, MVCNN has significantly advanced the state-of-the-art of 3D shape classification. The method renders a 3D shape to RGB or

depth images from different viewpoints, uses advanced CNN architecture to extract features for each view, and then aggregates the multi-view features to form the final feature representation based on max or average pooling. Albeit being simple yet effective, its best performance is achieved only when *all* views are used. In some practical scenarios, such as robot-operated active recognition, it is desirable to achieve object recognition with as-few-as-possible views, to minimize the robot movement cost.

Our key observation is that human is able to recognize a 3D shape without processing all views. Given the first view observation of a 3D shape, a human tends to first form hypotheses about which categories the shape may fall in, and then switch to the next viewpoint purposely, by moving himself around the shape or rotating the shape, to quick narrow down the uncertainty and refine hypotheses. This process is repeated for a few times until sufficient evidence are collected to minimize the uncertainty. There are two most prominent characters of the above procedure. First, the informative view is quite sparsely selected, far from being exhaustive. Second, both the next view selection and the predication making are deduced from the combined information from the previous observations over time.

Following the above intuition and drawing inspiration from visual attention model based on recurrent neural networks (RNN) [18], [19], we present *View-Enhanced Recurrent Attention Model (VERAM)* capable of automatically selecting an as short as possible sequence of views to classify 3D shapes. As shown in Fig. 1, our model is formulated as a goal-directed agent interacting with a 3D shape by using a camera sensor (Fig. 1(a)). At each time step, the agent actively selects the next viewpoint l for the camera sensor, according to the current internal state h . Next, the camera sensor captures the 3D shape from the new viewpoint to

- Corresponding author: Kai Xu (kevin.kai.xu@gmail.com).
- Songle Chen is with Jiangsu High Technology Research Key Laboratory for Wireless Sensor Networks, Nanjing University of Posts and Telecommunications. Email: chensongle@njupt.edu.cn.
- Lintao Zheng is with School of Computer, National University of Defense Technology. Email: lintaozheng1991@gmail.com.
- Yan Zhang is with Department of Computer Science and Technology, Nanjing University. Email: zhangyannju@nju.edu.cn.
- Zhixin Sun is with Jiangsu High Technology Research Key Laboratory for Wireless Sensor Networks, Nanjing University of Posts and Telecommunications. Email: sunzx@njupt.edu.cn.
- Kai Xu is with School of Computer, National University of Defense Technology. Email: kevin.kai.xu@gmail.com.

Manuscript received Sep. 27, 2017; revised May 26, 2018.

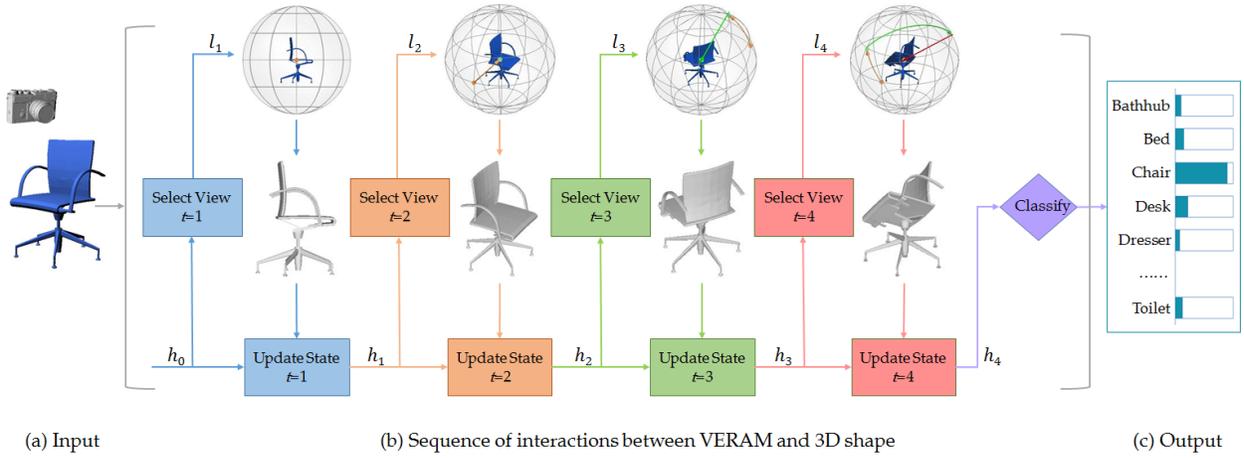


Fig. 1: VERAM is a view-enhanced recurrent attention model capable of adaptively selecting a sequence of views to classify 3D shapes. (a) The input is an unknown 3D shape to be rendered with a virtual camera. (b) The sequence of interactions between VERAM and the 3D shape. At each time step, VERAM actively selects the next view l for rendering according to current internal state h , which is then updated by the new observation. (c) Based on the aggregated information over time steps, VERAM makes a classification and outputs the category probabilities of the input 3D shape.

obtain a 2D image and pass it to the agent. The agent then extracts features for the new input image and updates its internal state. Such process is repeated for a few steps (4 time steps in Fig. 1(b)). Finally, based on the aggregated information over time steps, the agent emits the prediction as the probabilities of shape categories (Fig. 1(c)).

On the technical side, our method has two advantages over MVCNN. First, when pooling across all views, MVCNN discards the location information of each view. However, recent works reveal that viewpoint location plays an import role in enhancing the performance of the task of 3D shape classification [16], [17]. Second, processing all views involves high computational cost both at training and testing, although not every view is essential to recognition.

There have been a number of works using RNN-based visual attention model for image classification [18], [19], [20], [21], image captioning [22], and even for 3D object retrieval [23]. However, a common issue with these RNN-based attention models is the unbalanced training of the subnetworks corresponding to next view estimation and shape classification. The classification subnetwork is usually easier to train than the view estimation one. Consequently, the model training can be easily trapped in a local minima: the classification subnetwork is overfitted while the view estimation one poorly trained, which in turn leads to degraded classification accuracy. In this work, we propose three key technical contributions with VERAM, aiming to enhance the view learning and achieve a significant performance boost to multi-view 3D shape classification.

- We introduce three schemes to improve the information flow of gradient backpropagation from the view estimation subnetwork to the hidden units, achieving a balance with the training of the classification subnetwork. This overcomes the issue that the estimated views may get stuck at the boundary in the view parameterization space, which is usually encountered in previous works [19].

- During the reinforcement learning of our attention model, we integrate the classification confidence of the current view into the gradient computation of the reward against the view. This leads to a highly efficient guidance to the next view estimation.
- In VERAM, a novel loss function is proposed with a regularization term that enforces the estimated view to be distant to any of the previous ones so as to avoid view duplication.

The hybrid architecture of VERAM is trained for the subnetworks of shape classification and view estimation jointly, with the former using SGD [24] and the latter taking REINFORCE [25]. We empirically evaluate our model on the ModelNet benchmark [26]. Taking rendered gray-scale image as input and AlexNet as CNN architecture, without applying any data augmentation or network ensemble strategy, VERAM with 9 views achieves average instance-level and class-level accuracy of 95.5% and 95.3% on ModelNet10, 93.7% and 92.1% on ModelNet40. The high accuracy and efficiency make our model scalable to large datasets and applicable to many online applications.

2 RELATED WORK

3D shape classification via hand-craft descriptors. There is a long history of work in 3D shape analysis and a large variety of hand-craft shape descriptors have been presented. The representative view-based descriptors cover Light Field descriptor [1], Elevation descriptor [2], Aspect Graph based descriptor [3], [4] and DFT/DTW panoramic descriptor [5], etc. Popular shape descriptors include Shape Histogram descriptor [6], Spherical Harmonic descriptor [7], 3D SURF [8], Heat Kernel Signatures [9], etc. These descriptors are largely hand-engineered and usually do not have enough generalization ability to adapt to the diversity of numerous 3D shapes in different categories. As a result, their performance has an obvious gap compared with the current dominant methods based on deep learning technology, which

have achieved state-of-the-art performance in many tasks of computer graphics and computer vision. Our method falls into the deep learning class.

3D shape classification via deep CNN. A number of methods based on deep CNN have achieved state-of-the-art performance in 3D shape classification on public benchmarks. There are two categories. Shape-based methods [13], [26], [27], [28], [29], [30], [31], [32], [33] perform convolutions with 3D filters on the voxels or point clouds in continuous 3D space, and the volumetric representation makes them have the ability of exploiting complete structure information. View-based methods [12], [13], [14], [16], [17], [34], [35], [36], [37] first render the 3D shape into 2D images from different viewpoints, and then apply 2D filters to carry out convolution for each view.

Compared with shape-based methods, the advantage of view-based methods is that the massive image databases can be used to pre-train the deep neural network and advanced network architectures succeeded in image recognition tasks can be employed. Partly for these reasons, to data, view-based methods shows better or comparable performance to shape-based methods. Moreover, convolution on 2D images is more efficient than on 3D volumetric space. View-based methods naturally need to fuse clues from different 2D views, and max or average pooling is the most common strategy to perform the task [12], [13], [14], which lacks a view selection mechanism. Our method is view-based and also employs deep CNN to extract the descriptors for the rendered views. However, we adopt RNN-based visual attention model to learn attention policy of adaptively selecting a few number discriminative views, which is more effective and efficient.

3D shape recognition via active view selection. Our method fits into the realm of active recognition. Indeed, active recognition through next view planning has been studied for quite a long time in computer vision [38]. For 3D object recognition and pose estimation, next-best-view selection based on information rich model was proposed in [39]. In each step, the next-best-view is selected as the voxels of which has the highest number of matches that have not been detected before. In 3DShapeNets framework [26], next-best-view for 2.5D recognition is selected according to which can maximize the mutual information to reduce the potential uncertainty. In contrast to these local optimum next-best-view selection methods, an approximately global optimum approach was proposed by using undirected graph search [16]. However, the next-best-view selection is isolated from the neural network, and it is not a completely global optimum method. By contrast, the next-best-view predication of our method is embedded in deep RNN and is a global optimum approach. Recently, MV-RNN [23] combines RNN-based visual attention model with MVCNN [12] for 3D object retrieval. The view confidence and view location constrains are implicitly handled in the layer of feature representation. In contrast, VERAM does not combine MVCNN, and explicitly integrates view confidence and view location constrains into reward gradient and classification loss.

Visual attention model based on RNN. We draw inspiration from recent approaches that used RNN-based visual attention model to learn task-specific policies in various applications. The vision tasks include image classification

[18], [19], [20], [21], image caption generation [22], action with its boundary detection [40], and 3D object retrieval [23]. The attention model is also used in non-visual task, such as learning policies for a Neural Turing machine [41]. Our method builds on these directions and learns policies addressing the task for 3D shape classification. It extends the RNN-based visual attention models to be more robust to the common issue of the unbalanced training of the subnetworks, and provides a paradigm of how RNN-based attention model to support learning with view confidence and view location constrains integrally.

3 METHOD

The proposed VERAM is a RNN-based visual attention model for 3D shape classification, and is formulated as a goal-directed agent interacting with a 3D shape. A graphical representation of VERAM is shown in Fig. 2.

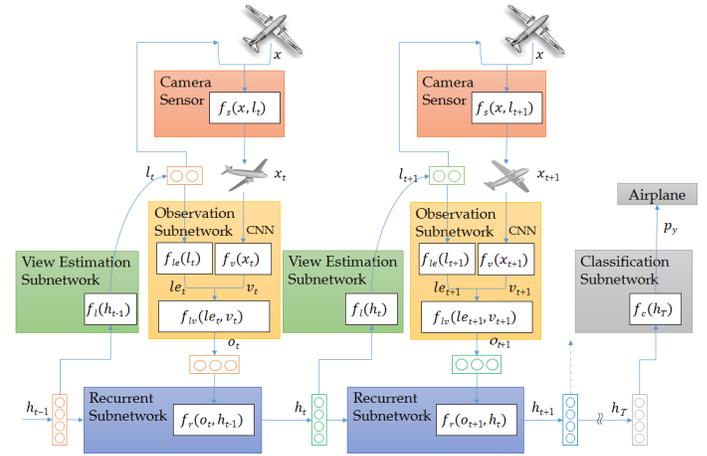


Fig. 2: Architecture of VERAM. Its sub-components include a virtual camera sensor, observation subnetwork, recurrent subnetwork, view estimation subnetwork and classification subnetwork.

The model sequentially processes a 3D shape x within T time steps. At each step t , based on the internal state h_{t-1} , the model actively selects a viewpoint l_t and obtains an observation x_t of 2D image rendered by the camera sensor from l_t , and then, the model uses x_t with l_t to updates its internal state. This process is repeated until the predication is emitted at step T . The architecture of the model will be described in subsection 3.1 and later in subsection 3.2, we explain how to use a combination of SGD and REINFORCE to train the model in end-to-end fashion.

3.1 Architecture

The architecture of VERAM can be broken down into a number of sub-components including camera sensor, observation subnetwork, recurrent subnetwork, view estimation subnetwork and classification subnetwork. Each component maps the input into a matrix or vector output.

3.1.1 Camera sensor

As shown in Fig. 2, the input to camera sensor is a 3D shape x and viewpoint location l_t , the output of the camera sensor is the 2D image x_t of the rendered view. 3D shape

x represented as polygon mesh is located at the center of the viewing sphere, as shown in Fig. 3 (a). The camera sensor can move on the surface of the viewing sphere and its location is indicated by the latitude and longitude. The camera sensor always points towards the centroid of the shape, and its upright vector is the tangent line of the latitude along clockwise direction. Phong reflection model is used to render the 3D shape into 2D images (224×224). Under a perspective projection, the pixel color is determined by the reflected intensity of the polygon vertices. Example rendered images are shown in Fig. 3 (b).

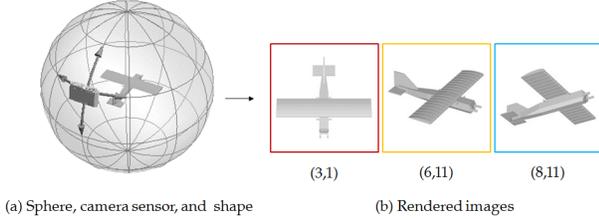


Fig. 3: 3D shape located in the center of the viewing sphere is rendered into 2D images by the camera sensor.

Data preparation is necessary for efficient training, and we sample discrete views at every 30 degrees both in latitude and longitude. If all views of a shape are rendered, they can be arranged in 12×12 grid, as shown in Fig. 4. The ordinal number 1 to 12 in the grid along vertical and horizontal direction is corresponding to the latitude and longitude of the sampled viewpoint location, which defines the view parameterization space for VERAM.

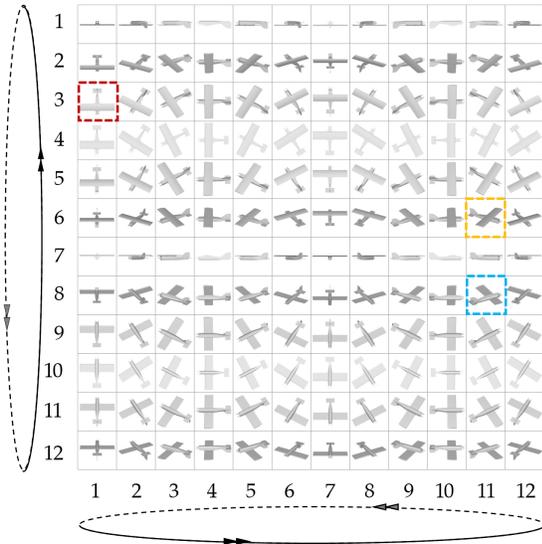


Fig. 4: Sample space of the viewpoint locations on viewing sphere is mapped to 12×12 grid for camera sensor, which defines the view parameterization space for VERAM.

There are two advantages of our sample strategy. First, the entire sphere is uniformly sampled both on latitude and longitude. Second, along the horizontal or vertical direction, the 12 images are consecutively connected and form a circle. As a result, it provides a contiguous space for the agent to deploy the camera sensor in an absolute or relative coordinate. In the following sections, the absolute viewpoint location in the view parameterization space is represented as $l = (r, c)$, and both r, c are in the range of $[1, 12]$. The

magnified images of (3, 1), (6, 11) and (8, 11) in Fig. 4 are shown in Fig. 3 (b). The function of the camera sensor can be formulated as

$$x_t = f_s(x, l_t). \quad (1)$$

3.1.2 Observation subnetwork

The job of the observation subnetwork is to encode the information about both where the observation is taken as well as what has been seen. As Fig. 2 illustrates, at each time step t , the observation subnetwork takes the rendered image x_t and location tuple $l_t = (r_t, c_t)$ as input and output a vector o_t .

We use $v_t = f_v(x_t; \theta_v)$ to denote the output v_t from function $f_v(\cdot)$ that takes image x_t as input and is parameterized by weights θ_v . $f_v(\cdot)$ typically maps with a sequence of convolutional, pooling, and fully connected layers, and the output of which is high level features v_t . Advanced network architectures succeeded in image recognition task can be used for $f_v(\cdot)$, such as AlexNet [10], VGG-16 [11], ResNet [42], etc. The location tuple l_t is mapped into embedding l_e by a fully connected layer $f_{l_e}(l_t; \theta_{l_e})$. In our practice, a fully connected layer is corresponding to a rectified linear unit $ReLU(Wx + b)$.

We concatenate the low bandwidth location l_e with the high bandwidth view information v_t by a fully connected layer $f_{lv}([l_e; v_t]; \theta_{lv})$ and output the final observation feature vector o_t . The observation subnetwork f_o can be represented as

$$o_t = f_o(l_t, x_t; \theta_o) = f_{lv}([f_{l_e}(l_t; \theta_{l_e}) f_v(x_t; \theta_v)]; \theta_{lv}), \quad (2)$$

where $\theta_o = [\theta_{l_e}, \theta_v, \theta_{lv}]$, corresponding to the whole parameters of the observation subnetwork.

3.1.3 Recurrent subnetwork

The agent maintains an internal state which encodes the agent's knowledge of the environment and summarizes information extracted from the history of past observations. It is instrumental to deciding how to act and where to deploy the camera sensor. In VERAM, this internal state is formed by the hidden units h_t of the recurrent neural network and updated over each time step with the feature vector o_t from the observation subnetwork, as shown in Fig. 2. The recurrent subnetwork is defined as

$$h_t = f_r(o_t, h_{t-1}; \theta_r). \quad (3)$$

linear mapping can be used for its efficiency, but Long-Short-Term Memory units (LSTM) [43] has the ability to learn long-range dependencies and stable dynamics.

3.1.4 View estimation subnetwork

The view estimation subnetwork acts as a controller that directs attention based on the current internal state. In Fig. 2, the view estimation subnetwork takes the hidden units h_{t-1} of recurrent subnetwork as input, and outputs l_t to make a prediction on where to deploy the camera sensor to render the next view. The view estimation subnetwork consists of $f_l(h_{t-1}; \theta_l)$ and $f_g(u_t; \theta_g)$. f_l maps the hidden units h_{t-1} into a two-dimensional coordinate tuple u_t , formally defined as

$$u_t = f_l(h_{t-1}; \theta_l). \quad (4)$$

f_l is usually implemented by a *Linear* layer followed by a specific transfer function. In the testing phase, u_t is directly used as the next viewpoint location l_t . The procedure of obtaining the image x_t from location l_t is non-differentiable, so in the training phase, a stochastic module $f_g(u_t; \theta_g)$ needs to be used. It samples l_t stochastically from a Gaussian distribution with a mean u_t and a fixed variance δ for reinforcement learning, defined as

$$l_t = f_g(u_t; \theta_g) = N(u_t, \delta^2). \quad (5)$$

The detail of view estimation subnetwork will be discussed in subsection 3.2.

3.1.5 Classification subnetwork

The classification subnetwork makes a classification and outputs the category probabilities y of the input 3D shape x based on the final internal state h_T , which integrates the information of the interaction history between the agent and input 3D shape. In VERAM, the classification subnetwork f_c has a fully connected layer followed by *LogSoftMax* output layer, namely

$$P(y|x) = f_c(h_T; \theta_c). \quad (6)$$

3.2 Learning

Given the category label y of shape x , we can formulate learning as a supervised classification problem. However, because the architecture of VERAM is hybrid, training it involves challenges of handling the non-differentiable component, of keeping balance between learning subnetworks, while struggling with the overfitting problem caused by millions of parameters. In this section, we will describe how VERAM combines SGD [24] and REINFORCE [25] to solve these problems.

3.2.1 REINFORCE-based learning

In subsection 3.1.1, the camera sensor is formulated as f_s to map 3D shape x with l_t to 2D image x_t . However, f_s is not a continuous function. As shown in Fig. 5, the gradient displayed as green arrow lines from observation subnetwork cannot back propagate to the view estimation subnetwork via the camera sensor. As a result, the view estimation subnetwork cannot be trained with standard back propagation.

REINFORCE [25] is adopted to solve this problem. Given a space of action sequences A , $p(a)$ is a distribution over $a \in A$ and parameterized by θ , we wish to learn network parameters θ that maximize the expected reward of action sequences. The gradient of the objective is

$$\nabla J(\theta) = \sum_{a \in A} p_\theta(a) \nabla \log p_\theta(a) r(a). \quad (7)$$

Here $r(a)$ is a reward assigned to each possible action sequence. Obviously, due to the high-dimensional space, this is a non-trivial optimization problem. REINFORCE addresses this by learning network parameters using Monte Carlo sampling. After running an agent's current policy π_θ in its environment and obtaining K interaction sequences of length T , the approximation to the gradient equation is

$$\nabla J(\theta) \approx 1/K \sum_{i=1}^K \sum_{t=1}^T \nabla \log \pi_\theta(a_t^i | s_{1:t}; \theta) (R_t^i - b_t). \quad (8)$$

Here, policy π_θ maps the history of past interactions with the environment $s_{1:t}$ to a distribution over actions for the current time step t . R_t is the cumulative future reward from the current time step to time step T , b_t is a baseline reward to reduce the variance of the gradient estimation. For a detail RNN-based REINFORCE, we refer to [25] and [44].

In our case, a_t is to predict the next location u_t , $s_{1:t}$ is summarized in the state of the hidden units h_{t-1} . $R = 1$ if the shape is classified correctly after T steps and 0 otherwise. Policy π is a fully connected hidden layer that maps h_{t-1} to u_t as formally defined in formula (4), and θ is corresponding to θ_l . Gaussian distribution is adopt to implement the reinforcement algorithm. Assuming f is the density function of the Gaussian distribution N as defined in formula (5), l_t is the sampled location, the gradient of R w.r.t. u_t is

$$\frac{dR}{du_t} = (R - b) \times \frac{d \ln(f(l_t, u_t))}{du_t}. \quad (9)$$

This gradient can easily back propagate to θ_l , so $\nabla J(\theta)$ can then be calculated. The flow direction of this gradient is shown as red arrow lines in Fig 5. REINFORCE learns model parameters according to this approximate gradient. It increases the log-probability of an action with a larger than expected cumulative reward, and decreases the probability if the obtained cumulative reward is smaller.

3.2.2 Enhancing the information flow of gradient

As shown in Fig. 5, observation subnetwork, recurrent subnetwork, and classification subnetwork are with standard deterministic neural network connections, and can be trained directly by back propagating gradients from the classification loss, namely, by SGD. By contrast, view estimation subnetwork containing stochastic module needs to be trained using REINFORCE described in subsection 3.2.1.

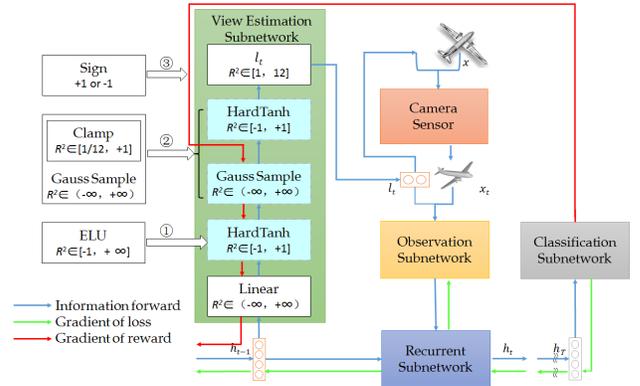


Fig. 5: Based on the visual attention model presented in [19], three schemes for enhancing information flow of gradient backpropagation for the view estimation subnetwork.

Giving deep insight into the architecture as shown in Fig. 2 and Fig. 5, we can find that there exists two ways starting from the hidden units h_{t-1} and ending at classification subnetwork. Analogous to other RNN-based visual attention models, VERAM is essentially a parallel neural network, and suffers from the common issue of the unbalanced training of the subnetworks corresponding to next view estimation and shape classification. Specifically,

as pointed in [21], the full model does not have enough time to learn a good attention policy before the classification subnetwork overfits to the data. In our practice, we find regardless of different initial parameters and different 3D shapes, the estimated view locations l_t tend to get stuck at the boundary of the view parameterization space, namely 1 or 12. Based on the visual attention model presented in [19], we propose three critical schemes to solve this issue, as shown in Fig. 5. It can be summarized as follows:

1) As mentioned in subsection 3.1.4, f_l is implemented by a *Linear* layer followed by a specific transfer function to force the estimated location u_t into the target range. *HardTanh* as well as *Sigmoid* are commonly used for their continuous interval can be easily mapped to view parameterization space. However, we find the problem of gradient disappearance caused by these transfer functions [45] is very serious, which means in reinforcement learning, the gradient from reward cannot back propagate to hidden units h . VERAM adopts *Eula* [46] as this specific transfer function, which can effectively alleviate the vanishing gradient problem via the identity for positive values.

2) Similar to f_l , we also need a function for f_g to force the output of *Gauss sample* to fall into the view parameterization space. As shown in Fig. 5, the second *HardTanh* operation is performed in [19] for this purpose. However, the gradient of reward bypasses it and directly applies to *Gauss sample* function. If the output l_t of *Gauss sample* is out of the range, the gradient calculated by formula (9) will be not right for l_t is not the actual viewpoint location. To solve this problem, VERAM performs *Clamp* in *Gauss sample* itself to simulate the output of *Gauss sample* is always in the range of view parameterization space.

3) In the backward process, formula (9) indicates that if the shape is misclassified, the learning process will encourage u_t move away from l_t . In the boundary, this also needs to be refreshed with sign as

$$sign = \begin{cases} +1, & \text{if } l_t \text{ is not boundary} \\ -1, & \text{if } u_t < l_t, l_t = 1/12, \text{left boundary} \\ -1, & \text{if } u_t > l_t, l_t = 1, \text{right boundary} \end{cases} . \quad (10)$$

Besides, if the shape is classified correctly, sign always sets to +1. Sign will multiplies to the result of formula (9). It means, in boundary, we still need to move u_t to l_t if u_t is not in the range of $[1/12, 1]$, although the shape is misclassified.

These three critical schemes enhance the information flow of gradient backpropagation from the view estimation subnetwork to the hidden units, and ensure each module be sequentially trained without break. As a result, it can effectively overcome the issue that the estimated views getting stuck at the boundary of the view parameterization space.

3.2.3 Learning with view confidence

The schemes for enhancing the information flow of gradient provides a basis for keeping a balance between the subnetworks. However, the classification subnetwork is still easier to train than the view estimation one, making the learned attention policy be easily trapped in local optimization. In this subsection, we propose a method of learning with view confidence for REINFORCE to solve this problem.

As mentioned in subsection 3.1.2, advanced network architectures succeeded in image recognition tasks can be used for $f_v(\cdot)$ to extract features v_t for image x_t . Besides, for all training shapes, we also extract the confidence c_t of image x_t . Concretely, first, we extract features v_t for each image x_t of all views of the training shapes. Second, each image x_t is also labeled with the same category y_t of the 3D shape. By this means, for each 3D shape in the training set, we can obtain 144 pairs of (v_t, y_t) . The collection of (v_t, y_t) of all the 3D shapes in the training set will be taken as the input to the following simple network

$$P(y_t|x_t) = \text{LogSoftMax}(\text{Linear}(v_t, \#categories)). \quad (11)$$

We use the negative log likelihood criterion to train this network. After convergence, the output probability of *LogSoftMax* corresponding to its category is extracted as the confidence c_t of image x_t . Fig. 6 presents the confidence of each 2D image shown in Fig. 4. The whiter the viewpoint, the more confidence the image has. We can see the image in (3,1) almost has no confidence for its own category *airplane*. According to the output of *LogSoftMax*, the highest category probability of this image is *monitor*.

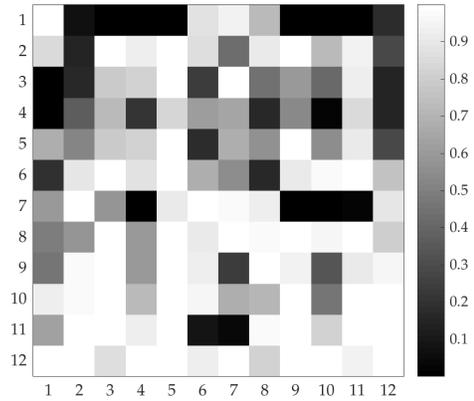


Fig. 6: Confidence, or category probability of each 2D image shown in Fig. 4.

If a shape is classified correctly, according to formula (9), reinforcement learning will encourage all u_{t+1} to move to l_t . On the contrary, it will encourage all u_{t+1} to move away from l_t . Intuitively, if a shape is classified correctly, but in step t , the confidence c_t of image x_t located at l_t is not high, we should weaken its effects of encouraging u_{t+1} to move to l_t . On the other hand, if a shape is misclassified, but the confidence c_t of image x_t located at l_t is high, we should weaken its effects of encouraging u_{t+1} to move away from l_t . Based on this inspiration, combining formula (10), we refresh formula (9) as

$$\frac{dR}{du_t} = sign \begin{cases} (R - b) \times \frac{d \ln(f(l_t, u_t))}{du_t} \times c_t \\ (R - b) \times \frac{d \ln(f(l_t, u_t))}{du_t} \times (1 - c_t) \end{cases} . \quad (12)$$

If a shape is classified correctly, the first line of formula (12) is used to calculate the gradient, otherwise, the second line of formula (12) is used. VERAM incorporates view confidence into reinforcement learning, which leads to a highly efficient guidance to the next view estimation, and

can effectively prevent the view estimation subnetwork from trapping into local optimization.

Note that confidence extraction is the preprocess procedure. The trained image classification network is only used to extract the confidence of image, and it will not be used anymore.

3.2.4 Learning with view location constrains

In practice, we find the visited view of each time step may overlap. To solve this problem and make better use of the complementary capacities of different views. A novel regular term for view location constrains is supplemented to loss function and learned at the same time to make the distribution of the visited view locations much more reasonable and mutually complementary.

The regular term for the view location constrains is depend on the prior knowledge of the applications. For 3D shape classification, there are total 144 view locations for the camera sensor, but VERAM only needs a few time steps T to form the judgment, so the visited views at least should be separated from each other. This weak regular term is adopted by VERAM.

Particularly, after the agent has moved through T steps, we add pairwise distance layer for each two view locations l_i and l_j the agent has visited. We train this part of the model with the loss of *HingeEmbeddingCriterion* to force the visited views to separate from each other. The loss is formally defined as

$$Loss(l_i, l_j) = \max(0, 1/12 - PairDistance(l_i, l_j)). \quad (13)$$

By integrating this loss function into the learning framework of recurrent attention model, VERAM can explicitly circumvents view duplication and the performance is further improved.

4 EXPERIMENTS AND DISCUSSIONS

4.1 Dataset, criteria and implementation details

Dataset. We evaluate VERAM along with current state-of-the-arts on ModelNet10 and ModelNet40 [26]. ModelNet10 contains 10 categories with 4899 unique 3D shapes, and ModelNet40 contains 40 categories with 12311 unique 3D shapes. The training set and testing set have been split on the website.

Criteria. We report classification accuracy with two level criteria. Instance-level accuracy is the ratio of the number of shapes that are classified correctly to the number of the total shapes. Class-level accuracy is the average of instance-level accuracy among all categories. Class-level accuracy is more objective for there are big difference among different categories in the number of testing shapes, from 20 to 100, but instance-level is more intuitive for it directly reflects how many shapes are misclassified.

Implementation details. VERAM is implemented by Torch on the platform with NVIDIA GeForce TITAN X. It needs about 1500 epochs for training. The learning rate is set to 0.001 in the first 600 epochs, then decreases linearly to minimum 0.00001 at epoch 1200. Momentum is set to 0.9. Based on grid search, the fixed variance δ of Gaussian distribution for REINFORCE is set to 0.11.

CNN architecture and recurrent subnetwork are two main components that affect the performance of VERAM. CNN is used to encode the rendered image. For CNN architecture, AlexNet [10] and ResNet [42] are used in our experiments. For recurrent subnetwork, linear mapping and LSTM [43] are adopted to update the hidden state in each time step. We will give the detail settings in each part of the evaluation.

Theoretically, the parameters of CNN in the observation subnetwork can be fine-tuned with the recurrent network at the same time. However, the training progress will be very slow because the forward and backward propagation need to proceed at each time step. For efficient training, the parameters of AlexNet and ResNet pre-trained on ImageNet are fixed and without further fine-tuning process on ModelNet. However, we can select more previous layer of CNN to counteract this influence.

The *source code* as well as the trained model of VERAM will be released at our project page: www.kevinkaixu.net/projects/veram.html.

4.2 Comparison with state-of-the-arts

In this subsection, we will compare the performance of VERAM with state-of-the-art deep neural-based methods, especially with the view-based methods. The performance of VERAM is achieved by taking AlexNet as CNN architecture for fair comparison. LSTM is adopted for recurrent subnetwork.

Comparison with deep neural-based methods. The performance of VERAM is compared with state-of-the-art deep neural-based methods on 3D shape classification, as summarized in table 1. These methods can be roughly grouped into two categories: shape-based and view-based. Shape-based methods cover 3DShapeNets [26], VoxNet [27], SubVolSup [13], AniProbing [13], VRN & VRN-Ensemble [28], FPNN [29], PointNet [30], PointNet++ [31], O-CNN [32] and So-Net [33]. The view-based methods include MVCNN [12], DeepPano [34], PANORAMA-NN [35], PANORAMA-ENN [36], MVCNN-Alex [13], MVCNN-MultiRes [13], Pairwise [16], DomSetClust [14], RotationNet [17] and the proposed VERAM. Besides, FusionNet [48] exploits both volumetric representation and projective pixel representation.

Among shape-based methods, VRN-Ensemble [28] achieves the best instance-level accuracy, 97.14% on ModelNet10 and 95.54% on ModelNet40. This result is by summing predictions from separately trained five VRN models and one Voxception model. When comparing VERAM with the single VRN model, VERAM gets the better performance, 95.5% vs. 93.61% on ModelNet10, and 93.7% vs. 91.33% on ModelNet40.

Among view-based methods, DeepPano [34], PANORAMA-NN [35] and PANORAMA-ENN [36] are based on the panoramic image of 3D shape. PANORAMA-ENN has a clear advantage over VERAM, 96.85% vs. 95.5% on ModelNet10, 95.56% vs. 93.7% on ModelNet40. PANORAMA-ENN needs to obtain three panoramas from different principle axes and each panorama needs to extract SDM, NDM and magnitude of gradient image of NDM, while VERAM only takes grayscale image as input. Moreover, these three panoramic methods can be regarded

TABLE 1: Comparison of classification accuracy of methods based on deep neural networks on ModelNet10 & 40.

	Method	ModelNet10		ModelNet40	
		Inst.	Class	Inst.	Class
shape-based	3DShapeNets [26]	-	83.5	-	77.3
	VoxNet [27]	-	92.0	-	83.0
	SubVolSup [13]	-	-	89.2	86.0
	AniProbing [13]	-	-	89.9	85.6
	VRN [28]	93.61	-	91.33	-
	VRN-Ensemble [28]	97.14	-	95.54	-
	FPNN [29]	-	-	88.4	-
	PointNet [30]	-	-	89.2	86.2
	PointNet++ [31]	-	-	91.9	-
	O-CNN [32]	-	-	90.6	-
So-Net [33]	95.7	95.5	93.4	90.8	
view-based	DeepPano [34]	-	88.7	-	82.5
	PANORAMA-NN [35]	91.12	-	90.70	-
	PANORAMA-ENN [36]	96.85	-	95.56	-
	MVCNN [12]	-	-	-	90.1
	MVCNN-Alex [13]	-	-	92.0	89.7
	MVCNN-MultiRes [13]	-	-	93.8	91.4
	Pairwise [16]	94.0	-	92.0	-
	DomSetClust [14]	-	-	93.8	92.8
	RotationNet [17]	98.46	-	97.37	-
	VERAM	95.5	95.3	93.7	92.1
	mix FusionNet [48]	93.1	-	90.8	-

as based on continuous views, while VERAM and other view-based methods are based on discrete views.

Table 2 gives the detail comparison against state-of-the-art discrete view-based methods with their different processing strategies. Apparently, RotationNet, DomSetClust and MVCNN-MultiRes get the better performance. By augmenting the classification task with pose estimation, RotationNet gets instance-level accuracy 98.46% on ModelNet10 and 97.37% on ModelNet40, both are state-of-the-art performance. They are achieved by alerting 11 different camera system orientations. According to table 7 of [17], with AlexNet, On ModelNet 40, there are 8 of total 11 camera system orientations (except 2th, 3th, 4th) under which the performance of RotionalNet is less than 93.04%, which is inferior to VERAM (93.7%). Considering VERAM only uses a single camera system orientation, there is potential for VERAM to further reduce the gap.

DomSetClust ever gets instance-level accuracy 93.8% and class-level accuracy 92.8% on ModelNet40, and both better than the performance of VERAM. However, to obtain such performance, DomSetClust needs to take grayscale, depth and surface normals image as input and fine-tune the CNN network. When DomSetClust and VERAM both use grayscale image as input, VERAM has an clear advantage, 93.7% vs. 92.2% instance-level accuracy, and 92.1% vs. 91.5% class-level accuracy.

The comparison of VERAM and MVCNN-Alex is more fair for they both taking grayscale image as input, AlexNet as CNN, and with single resolution. VERAM has an obvious advantage over MVCNN-Alex, 93.7% vs. 92.0% instance-level accuracy and 92.1% vs. 89.7% class-level accuracy. The performances of VERAM is matchable with MVCNN-MultiRes, but the latter needs to implement two times of voxelization and three times of rendering.

Based on the above study, it can be concluded that under the equivalent conditions of render representation, resolution, CNN architecture and number of views, VERAM outperforms all state-of-the-art view-based methods.

Comparison with alternative view-selection methods. VERAM falls into the category of active view selection. Among the methods in table 1, both 3DShapeNets [26] and Pairwise [16] also adopt active view selection for 3D shape classification. Besides, MV-RNN [23] extends RNN-based visual attention model by integrating MVCNN [12] into the architecture for 3D object retrieval. To give a comprehensive comparison, we implemented MV-RNN for 3D shape classification and the input to MV-RNN is 2048 vector for each rendered image extracted from ResNet152. Both VERAM and MV-RNN take grayscale image as input, and 3DShapeNets takes depth image as input. In contrast, Pairwise exploits more information and takes both grayscale and depth image as input.

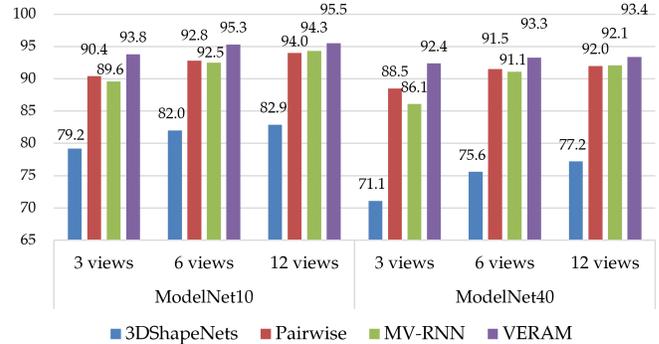


Fig. 7: Comparison of instance-level accuracy of 3DShapeNets [26], Pairwise [16], MV-RNN [23] and VERAM with different number of views on ModelNet10 and ModelNet40.

3DShapeNets provides the baseline performance on both datasets among all methods. With the deep neural networks pre-trained on ImageNet, Pairwise makes a great progress on the performance. With 12 views, it gets 94.0% and 92.0% instance-level accuracy on ModelNet10 and ModelNet40 respectively. MV-RNN adopts CNN2 of max pooling to perform classification and the accuracy increases with the number of views. The performance of MV-RNN outperforms that of Pairwise with 12 views. VERAM outperforms 3DShapeNets, Pairwise and MV-RNN on both datasets in every view. Although VERAM also uses the pre-trained model on ImageNet. However, the next view selection of VERAM is embedded seamlessly in RNN with view confidence and view location constrains learning.

4.3 Evaluation of the enhancement of VERAM

In this subsection, we will evaluate the three key technical components of VERAM on conducting to improve the performance step by step. ResNet is used to encode the rendered image for it is more compact. We use linear mapping for recurrent network instead of LSTM to stress the ability of three key components to select more discriminative views to enhance the view learning. In the end of this subsection, we will present the results when using LSTM as recurrent

TABLE 2: Detailed comparison against state-of-the-art discrete view-based methods on ModelNet10 & ModelNet40.

Method	Input	Reso- lution	CNN	Fine tune	ModelNet10				ModelNet40			
					View	Inst.	View	Class	View	Inst.	View	Class
MVCNN [12]	Gray	1	VGG-M	✓	-	-	-	-	-	-	80	90.1
MVCNN-Alex [13]	Gray	1	AlexNet	✓	-	-	-	-	20	92.0	20	89.7
MVCNN-MulRes [13]	Gray	3	AlexNet	✓	-	-	-	-	20	93.8	20	91.4
Pairwise [16]	Gray+Depth	1	VGG-M	✓	12	94.0	-	-	12	92.0	-	-
DomSetClust [14]	Gray	1	VGG-M	×	-	-	-	-	12	91.9	12	90.4
DomSetClust [14]	Gray+Depth+Surf	1	VGG-M	×	-	-	-	-	12	93.3	12	92.1
DomSetClust [14]	Gray	1	VGG-M	✓	-	-	-	-	12	92.2	12	91.5
DomSetClust [14]	Gray+Depth+Surf	1	VGG-M	✓	-	-	-	-	12	93.8	12	92.8
RotationNet [17]	Gray	1	AlexNet	✓	20	98.46	-	-	20	97.37	-	-
VERAM	Gray	1	AlexNet	×	9	95.5	9	95.3	9	93.7	9	92.1
VERAM	Gray	1	ResNet	×	9	96.3	9	96.1	9	93.2	9	91.5

network. We trained 5 models for each network setting with the specified super parameters, and use the average class-level accuracy for comparison.

Enhancing the information flow of gradient. RNN-based visual attention model suffers from the problem of unbalanced training of the subnetworks, and the estimated view locations tend to get stuck at the boundary in the view parameterization space. Fig. 8 (left) shows the heat map of view location frequency of each time step by applying the visual attention model presented in [19], denoted as ClassicalRAM, to predict all *chairs* in the testing set of ModelNet10, time steps $T=4$. It can be seen that starting from the view located at (6, 4), the views of the next three steps almost all locate at the boundary (1, 1). To solve this issue, three schemes are proposed for VERAM as described in subsection 3.2.2 to enhance information flow of gradient backpropagation. For simplicity, we call it BoundaryRAM. Fig. 8 (right) shows the heat map of view location frequency of each time step by applying a trained BoundaryRAM model to the same *chair* dataset. There are significant difference among the heat maps of different time steps.

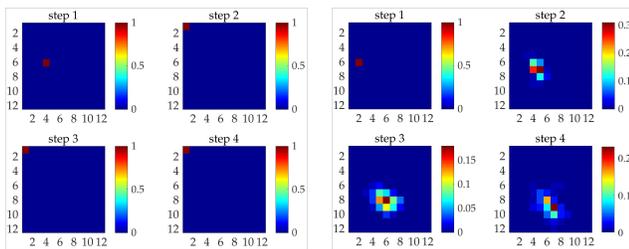


Fig. 8: Heat maps of view location frequency of predicating all *chairs* in the testing set of ModelNet10, time steps $T=4$. Left, by using ClassicalRAM [19]. Right, by using BoundaryRAM.

The comparison of the average class-level accuracy between ClassicalRAM [19] and BoundaryRAM on ModelNet10 and ModelNet40 is shown in Fig. 9. The horizontal axis is time steps T from 1 to 6. Note that T is a super parameter, and the accuracy of each T is the average value from 5 trained models with the same network setting. Leaving out $T=1$, i.e., $T=2 \sim 6$, the accuracy of BoundaryRAM goes up from min 1.80% to max 5.37% on ModelNet10, and from min 0.65% to max 2.76% on ModelNet40. From Fig. 8 and Fig. 9, it can be seen that the three critical schemes

presented in subsection 3.2.2 can effectively overcome the issue that the estimated views getting stuck at the boundary, and achieve a apparent performance enhancement.

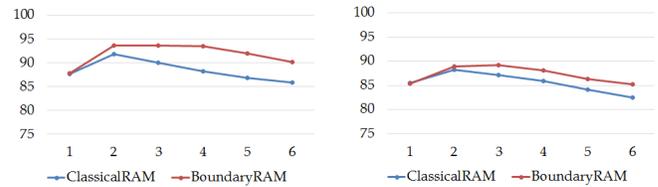


Fig. 9: Comparison of average class-level accuracy between ClassicalRAM [19] and BoundaryRAM. Left, on ModelNet10. Right, on ModelNet40.

Learning with view confidence. The classification subnetwork of VERAM is easier to train than the view estimation one, and learned attention policy can be easily trapped in local optimization. In practice, with the same network, we find the accuracy of each category from different trained BoundaryRAM model fluctuates widely. With the number of time steps increases, the model trained by BoundaryRAM will be more unstable. To alleviate this problem, based on BoundaryRAM, in subsection 3.2.3, we propose a method of learning with view confidence for REINFORCE to provide effective guidance to agent on where to deploy the model’s attention, here called ConFRAM. The comparison of the average class-level accuracy among ClassicalRAM [19], BoundaryRAM and ConFRAM on ModelNet10 and ModelNet40 is shown in Fig. 10.

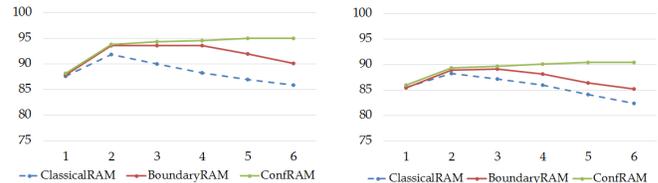


Fig. 10: Comparison of average class-level accuracy among ClassicalRAM [19], BoundaryRAM, and ConFRAM. Left, on ModelNet10. Right, on ModelNet40.

From Fig. 10, it can be seen that, the performance of ClassicalRAM drops after the number of time steps $T=2$, while BoundaryRAM drops after $T=4$ on ModelNet10 and after $T=3$ on ModelNet40. By contrast, ConFRAM can exploit more views and converge at about $T=6$ with signif-

icant higher performance. Compared with BoundaryRAM, leaving out $T=1$, i.e., $T=2 \sim 6$, the accuracy of ConfRAM goes up from min 0.20% to max 4.87% on ModelNet10, and from min 0.36% to max 5.24% on ModelNet40.

Learning with view location constrains. ConfRAM can achieve a stable and fairly good performance. However, the visited view location of each time step may overlap. In subsection 3.2.4, based on BoundaryRAM and ConfRAM, a weak regular term is adopted by VERAM to keep the visited views separated from each other, here called LocRAM, equally to the whole VERAM. The comparison of the average class-level accuracy among ClassicalRAM [19], BoundaryRAM, ConfRAM and LocRAM on ModelNet10 and ModelNet40 is shown in Fig. 11.

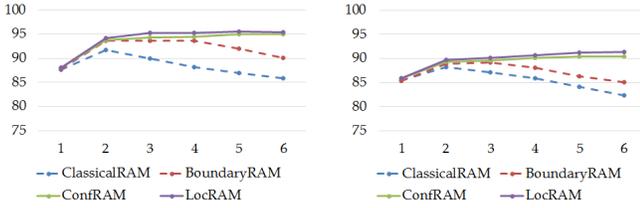


Fig. 11: Comparison of average class-level accuracy among ClassicalRAM [19], BoundaryRAM, ConfRAM and LocRAM. Left, on ModelNet10. Right, on ModelNet40.

Fig. 11 shows LocRAM learning with the weak regular term of view location constrains can obtain a clear performance improvement. Compared with ConfRAM, leaving out $T=1$, i.e., $T=2 \sim 6$, the accuracy of LocRAM goes up from min 0.40% to max 0.94% on ModelNet10, and from min 0.43% to max 0.91% on ModelNet40. It should be pointed out that the regular term of VERAM is only a weak constrain.

The above experiments of this subsection use linear mapping as recurrent subnetwork, and the previous step will have less impact when the interval between which and the last step T increases. When replacing linear mapping with LSTM, the comparison of the average class-level accuracy among ClassicalRAM [19], BoundaryRAM, ConfRAM and LocRAM on ModelNet10 and ModelNet40 is shown in Fig. 12. Since LSTM is effective at capturing long-term temporal dependencies, the already selected discriminative view prevents the performance from decreasing obviously. However, VERAM with three key components still achieves a significant performance boost over ClassicalRAM. When $T=6$, the accuracy is improved by 3.23% on ModelNet10 (95.54% vs. 92.31%), and by 3.12% on ModelNet40 (91.44% vs. 88.32%).

4.4 Affecting factors on the performance of VERAM

Effect of different recurrent subnetworks. The final predication is based on the hidden state of recurrent subnetwork which is updated over each time step. Linear mapping and LSTM are mostly used for this purpose. When using linear mapping as recurrent subnetwork, formula (3) is specified as $h_t = ReLU(Linear(h_{t-1}) + o_t)$. For LSTM, the total units is set to 1024 and each of which is composed of cell, input gate, output gate and forget gate. Fig. 13 gives

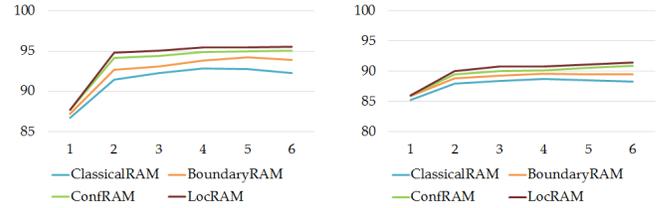


Fig. 12: Comparison of average class-level accuracy among ClassicalRAM [19], BoundaryRAM, ConfRAM and LocRAM when using LSTM as recurrent cells. Left, on ModelNet10. Right, on ModelNet40.

the comparison of the performance of VERAM with linear mapping and with LSTM on ModelNet40 when $T=3, 6, 9$. The left uses AlexNet as CNN while the right uses ResNet as CNN.

From Fig. 13, it can be seen that LSTM achieves the better performance than linear mapping. When using AlexNet to encode the visual representation, VERAM with LSTM has significant advantage over with linear mapping. For example, the instance-level accuracy increases 3.16% from 90.56% to 93.72% by LSTM when $T=9$. When using ResNet as CNN, LSTM still achieves the better performance than linear mapping. However, the improvement is slight. When $T=9$, it only gets instance-level accuracy 0.52% increment from 92.67% to 93.19%. Besides, for both AlexNet and ResNet, the performance of linear mapping when $T=9$ is slightly lower than when $T=6$, but the performance of LSTM is quite the opposite.

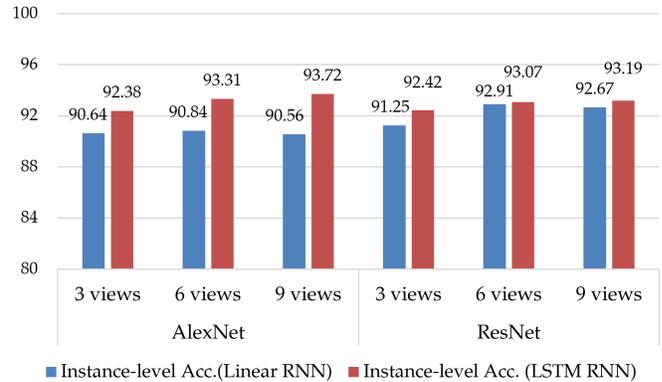


Fig. 13: Comparison of the best instance-level accuracy of VERAM with linear mapping and with LSTM on ModelNet40 when $T=3, 6, 9$. Left, using AlexNet as CNN. Right, using ResNet as CNN.

Effect of different CNN architectures. As described in subsection 3.1.2, VERAM needs to encode the visual representation for the rendered image at each time step. A number of advanced CNN architectures can be adopted for implementation. Among them, AlexNet is the first notably successful CNN architecture on ImageNet while ResNet is the recently proposed representing the advanced level, and they are compared in this subsection. For each 2D image x_t of 3D shape, we extract features $v_t \in R^{4096}$ from layer *fc6* of the AlexNet pre-trained on ImageNet. Analogously, we extract features $v_t \in R^{2048}$ from layer *flatten0_output* of the ResNet152 pre-trained on ImageNet 11K, which is available on MXNet [47]. For efficient training, we didn't

fine-tune AlexNet or ResNet model on ModelNet. Fig. 14 presents the performance of VERAM with AlexNet and with ResNet on ModelNet40 when $T=3, 6, 9$. The left uses linear mapping as recurrent network while the right uses LSTM.

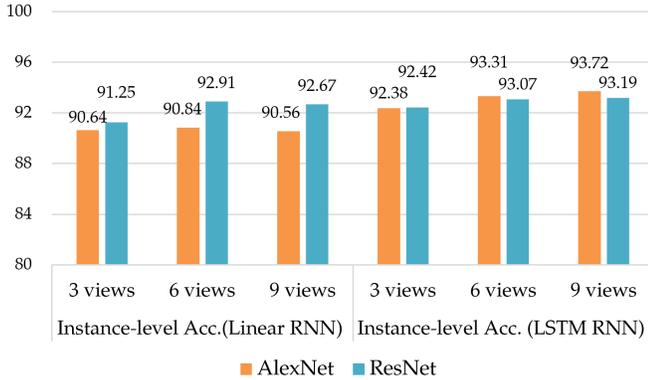


Fig. 14: Comparison of the best instance-level accuracy of VERAM with AlexNet and ResNet on ModelNet40 when $T=3, 6, 9$. Left, using linear mapping as RNN. Right, using LSTM as RNN.

Fig. 14 indicates that when using linear mapping as RNN, the performance of ResNet soars about 2 percent compared with that of AlexNet and has a clear advantage. However, when using LSTM as RNN, the performance of AlexNet and ResNet are comparable to each other. At $T=9$, the best instance-level accuracy of VERAM with AlexNet is 93.72%, and with ResNet, it is 93.19%. Theoretically, the feature extracted from ResNet has more capacity of discrimination than feature extracted from CNN in each time step are further merged by LSTM, and according to Fig. 13, the improvement of LSTM on VERAM with AlexNet is much more significant than on VERAM with ResNet. As a result, the performance gap between these two CNNs is reduced. We noticed the performance difference among different CNN architectures in recent work [17] on 3D shape recognitions is about 1%.

Effect of time steps T . Time steps T in VERAM is a super parameter and it determines how many images does VERAM need to render and observe before emit the classification predication. Fig. 15 shows the best instance-level and class-level accuracy VERAM achieved on ModelNet40 with different time steps T from 1 to 16. AlexNet is used to encode the visual representation and LSTM is adopt as recurrent subnetwork. For clarity, we only append T and its accuracy for 1, 3, 6, 9, 12, 16. In Fig. 15, the rate of accuracy growth is fast in the first few steps, then becomes slower and slower with the increase of time step T and converges at $T=9$, where VERAM achieves the best performance, 93.7% instance-level accuracy and 92.1% class-level accuracy. After that, the performance even slightly decreases. From Fig. 15 it can be concluded as follows:

- 1) The performance of VERAM can quickly converge within a few time steps. On ModelNet40, VERAM can get 92.38% instance-level accuracy with only 3 views, which is as high as 98.6% of the best performance (93.72%).
- 2) Increasing T would do little to improve the performance after the convergence. VERAM uses the same parameters and network for each time step, so that the total

number of parameters will not expand with the increase of T . Although increasing T means VERAM can obtain more information, but the larger T means VERAM needing to cope with more views with the same number of parameters.

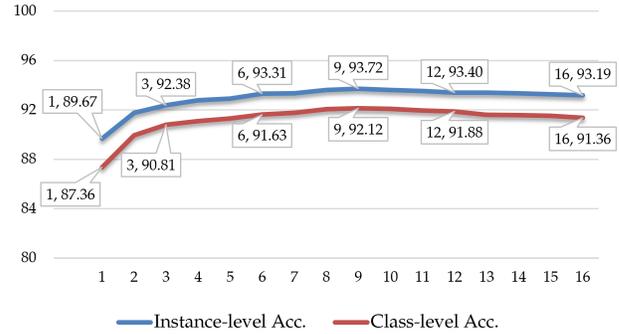


Fig. 15: Best instance-level and class-level accuracy of VERAM on ModelNet40 with time steps T from 1 to 16. AlexNet is used to encode the visual representation and LSTM is adopt as recurrent network.

Effect of shape alignment. As Pairwise [16] and RotationNet [17], VERAM needs a unified coordinate system to render shapes and each shape should be rendered from pre-defined viewpoints, so shape alignment is important to the proposed method. On the website of ModelNet [26], all shapes in ModelNet10 are manually cleaned and their orientations are well aligned. The shapes in ModelNet40 are also cleaned, but a small part of shapes are needed to align in the pre-processing stage. The reported performance of VERAM is based on such aligned orientation.

To quantify the effect of the alignment on the performance, we conducted an experiment with aligned ModelNet40 newly released on the website of ModelNet. Taking AlexNet as CNN and LSTM as RNN, we trained three VERAM models ($T=3, T=6$ and $T=9$) with the aligned training dataset and tested the classification on both aligned and unaligned (i.e., randomly rotated) shapes. The results are shown in Fig. 16. The instance-level accuracy of unaligned shapes only got 71.64%, 75.24% and 78.04% with 3, 6 and 9 views. Moreover, it only increases slightly with the increment number of views. These indicate that VERAM is rather sensitive to the pre-defined viewpoints. We notice that this limitation is generally encountered in several view-based methods. For instance, RotationNet [17] achieves 38% on ShapeNetCore55 dataset. In future, we plan to merge alignment network into VERAM to alleviate this problem.

Complexity. With the extracted features from AlexNet and using LSTM as RNN, the total number of parameters of VERAM is about 12.8M, which is lesser than that of VGG16 (14.7M) and ResNet152 (60M). In our single TITAN X GPU platform and in the testing phase, it takes less than 0.003 second to extract features for an image by AlexNet. Besides feature extraction, rendering and predication run very fast, and it only needs about 0.005 second to emit the decision. As shown in Fig. 15, the performance of VERAM is very stable with different T after $T > 3$. In summary, if we set $T=6$, which can already obtain high accuracy, it takes about 0.025 second to classify a 3D shape (not including the time spent on moving the camera).

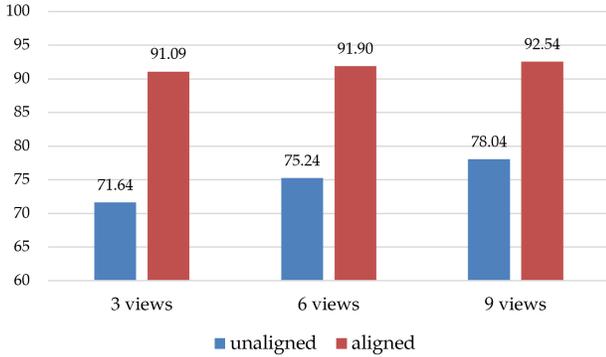


Fig. 16: Comparison of the instance-level accuracy of VERAM on unaligned and aligned test dataset of ModelNet40.

4.5 Detailed statistics

VERAM achieves the best average class-level accuracy 96.1% on ModelNet10 (with ResNet) and 92.1% on ModelNet40 (with AlexNet). The number of correctly classified shapes over the total number of shapes of each category with these two best models are reported in table 3.

On ModelNet10, *nightstand* and *table* are the two worst categories and they achieve 86.05% and 92% class-level accuracy respectively. 12 out of 86 *nightstands* in the testing set are misclassified. One each of them is misclassified as *desk* and *table*, and the other 10 *nightstands* are misclassified as *dresser*. A total of 8 *tables* are misclassified. One each of them is misclassified as *bathhub* and *dresser* and the others are misclassified as *desk*. Two misclassified shapes of *nightstand* and *table* are shown in the first row of Fig. 17.

TABLE 3: The number of correctly classified shapes over the total number of shapes of each category on ModelNet10 and ModelNet40 with the best class-level accuracy VERAM.

ModelNet10	bathtub	bed	chair	desk	dresser
	48/50	99/100	100/100	83/86	81/86
	monitor	nightstand	sofa	table	toilet
	100/100	74/86	97/100	92/100	100/100
	airplane	bathtub	bed	bench	bookshelf
100/100	47/50	100/100	16/20	98/100	
bottle	bowl	car	chair	cone	
93/100	17/20	100/100	98/100	19/20	
cup	curtain	desk	door	dresser	
12/20	19/20	81/86	19/20	78/86	
ModelNet40	flowerpot	glassbox	guitar	keyboard	lamp
	17/20	94/100	100/100	20/20	18/20
	laptop	mantel	monitor	nightstand	person
	20/20	96/100	100/100	62/86	20/20
	piano	plant	radio	rangehood	sink
	95/100	95/100	19/20	93/100	18/20
	Sofa	stairs	stool	table	tent
	97/100	18/20	18/20	79/100	19/20
	toilet	tvstand	vase	wardrobe	xbox
	100/100	92/100	89/100	18/20	16/20

On ModelNet40, the classification accuracy of *cup* is 60.0% and is the worst among all categories. *nightstand* gets 72.1% accuracy and holds the second worst position. Total 8 *cups* among 20 are misclassified. One each of them is misclassified as *bottle*, *bowl* and *wardrobe*, and the other 5 *cups* are misclassified as *vase*. The testing set of ModelNet40

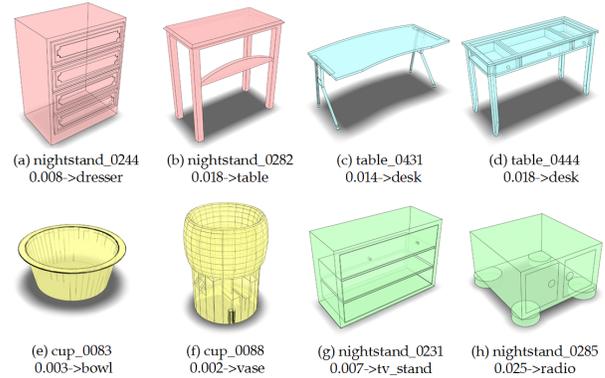


Fig. 17: First row shows the examples of misclassified shapes of *nightstand* and *table* on ModelNet10. Second row shows the examples of misclassified shapes of *cup* and *nightstand* on ModelNet40. Shape name, probability to its own category and misclassified category are denoted below the shape.

contains 86 *nightstands* and 24 are misclassified. Among them, 15 shapes are misclassified as *dresser* and the other 9 *nightstands* are spread to 5 categories including *bookshelf*, *glassbox*, *radio*, *table*, *tvstand*. Two misclassified shapes of *cup* and *nightstand* are shown in the second row of Fig. 17.

According to our intuition, the shapes shown in Fig. 17 are hard to be correctly classified by only using the visual representations. At least partly, the mistake due to a shape may have diversified functions. We suspect that more structure information in 3D space is needed to be fully exploited to meet such challenges.

5 CONCLUSION

We have presented VERAM, a recurrent attention model capable of actively selecting a sequence of views for highly accurate 3D shape classification. To address the problem commonly found in existing RNN-based attention models, i.e., the unbalanced training of the subnetworks, we propose to 1) enhance the information flow of gradient backpropagation for the view estimation subnetwork, 2) design a highly informative reward function for the reinforcement training of view estimation, and 3) formulate a novel loss function that explicitly circumvents view duplication.

When being applied to real scenarios, VERAM has several limitations which may spark future research:

Shape alignment. The experiment in subsection 4.4 shows that the result is sensitive to alignment and this issue is also reported in [17]. To address this problem, the basic strategy is by deepening the network and augmenting the training data to force the network to cover different viewpoint variants. Inspired by the approaches of learning transform parameters as [49], we think the more promising approach is to merge the alignment network into VERAM.

Time steps. VERAM uses the fixed time steps for prediction. To make the model capable of stopping observation once it has enough information, the reward of MV-RNN [23] contains the information gain and it terminates the process when entropy is less than the threshold, while [40] designed a subnetwork to learn a binary prediction indicator. Both of

them provide the insight to extend our model with varying instead of fixed time steps.

Inaccessible views and occlusion. In real scanning scenario, there are cases where some predicted views are physically inaccessible. Varying time steps is necessary and the proposed method also needs to be extended to learn the relevance of each view to the task as TAGM [50]. Occlusion also has severe adverse effect since the feature representation of each view is by convoluting the entire image. The key to this problem is to encode the part-based feature as in [23] to spot informative visible parts of the partially occluded 3d shape.

Moving cost and views passed through. The cost of moving scanner should be considered carefully in real scanning scenario. One feasible approach is to model the cost of moving scanner as the circle distance between adjacent selected views and add the cost to the reward for reinforce learning. Another issue is VERAM omits the continuous images obtained when moving the scanner from the current selected viewpoint to the next one. Although these views are not the best for next observation, but they also have the potential to help the classification to be more efficient.

Visual feature encoding. The pre-trained deep CNN is adopted by VERAM to extract visual features. For efficiency reason, such deep network is hard to fine-tuned with the view estimation subnetwork simultaneously, although what and where to observe are coupled for each other. It seems that visual attention model should exploit a more efficient network to learning what to observe.

ACKNOWLEDGMENTS

This work was supported in part by the National Natural Science Foundation of China (61373135, 61672299, 61702281, 61532003, 61572507 and 61622212), and the Post-doctoral Science Foundation of Jiangsu Province of China (No.1701046A).

REFERENCES

- [1] D. Y. Chen, X. P. Tian, Y. T. Shen, and M. Ouhyoung, "On visual similarity based 3D model retrieval," *Comput. Graph. Forum*, vol. 22, no. 3, pp. 223-232, 2003.
- [2] J. L. Shih, C. H. Lee, and J. T. Wang, "A new 3D model retrieval approach based on the elevation descriptor," *Pattern Recogn.*, vol. 40, no. 1, pp. 283-295, 2007.
- [3] C. M. Cyr and B. B. Kimia, "3D object recognition using shape similarity-based aspect graph," *Proc. IEEE Conf. Comput. Vis.*, pp. 254-261, 2001.
- [4] M. Ulrich, C. Wiedemann, and C. Steger, "Combining scale-space and similarity-based aspect graphs for fast 3d object recognition," *IEEE trans. on Pattern Analysis and Machine Intelligence*, vol. 34, no. 10, pp. 1902-1914, 2012.
- [5] P. Papadakis, I. Pratikakis, T. Theoharis, and S. Perantonis, "PANORAMA: A 3D shape descriptor based on panoramic views for unsupervised 3D object retrieval," *Int. J. Comput. Vision*, vol. 89, no. 2, pp. 177-192, 2010.
- [6] M. Ankerst, G. Kastentmiller, H. P. Kriegel, and T. Seidl, "Nearest neighbor classification in 3D protein databases," *Proc. the Seventh International Conference on Intelligent Systems for Molecular Biology*, pp. 34-43, 1999.
- [7] M. Kazhdan, T. Funkhouser, and S. Rusinkiewicz, "Rotation invariant spherical harmonic representation of 3d shape descriptors," *Proc. Symp. Geometry Process*, pp. 156-164, 2003.
- [8] J. Knopp, M. Prasad, G. Willems, R. Timofte and L. Van Gool, "Hough transform and 3D SURF for robust three dimensional classification," *Proc. 11th Eur. Conf. Comput. Vis.*, pp. 589-602, 2010.
- [9] A. M. Bronstein, M. M. Bronstein, L. J. Guibas, and M. Ovsjanikov, "Shape Google: Geometric Words and Expressions for Invariant Shape Retrieval," *ACM Trans. Graphics*, vol. 30, no. 1, pp. 1-20, 2011.
- [10] A. Krizhevsky, I. Sutskever, and G. Hinton, "ImageNet classification with deep convolutional neural networks," *Proc. Adv. Neural Inf. Process. Syst.*, pp. 1097-1105, 2012.
- [11] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv:1409.1556*, 2014.
- [12] H. Su, S. Maji, E. Kalogerakis, and E. G. Learned-Miller, "Multi-view convolutional neural networks for 3D shape recognition," *Proc. Int. Conf. Comput. Vision*, pp. 945-953, 2015.
- [13] C. R. Qi, H. Su, M. Niessner, A. Dai, M. Yan and L. J. Guibas, "Volumetric and multi-view CNNs for object classification on 3D data," *Proc. 2016 IEEE Conf. Computer Vision and Pattern Recognition*, pp. 5648-5656, 2016.
- [14] C. Wang, M. Pelillo, and K. Siddiqi, "Dominant Set Clustering and Pooling for Multi-View 3D Object Recognition," *Proc. 2017 British Machine Vision Conference (BMVC)*, 2017.
- [15] J. Deng, W. Dong, R. Socher, L. J. Li, K. Li and L. Fei-Fei, "Imagenet: A Large-Scale Hierarchical Image Database," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, pp. 248-255, 2009.
- [16] E. Johns, S. Leutenegger, and A. J. Davison, "Pairwise decomposition of image sequences for active multi-view recognition," *Proc. 2016 IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 3813-3822, 2016.
- [17] A. Kanazaki, Y. Matsushitay and Y. Nishida, "RotationNet: Joint Object Categorization and Pose Estimation Using Multiviews from Unsupervised Viewpoints," *Proc. 2018 IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 5010-5019, 2018.
- [18] J. Ba, V. Mnih, and K. Kavukcuoglu, "Multiple object recognition with visual attention," *arXiv:1412.7755*, 2014.
- [19] V. Mnih, N. Heess, A. Graves, and K. Kavukcuoglu, "Recurrent models of visual attention," *Proc. Adv. NIPS*, pp. 1-9, 2014.
- [20] P. Sermanet, A. Frome, and E. Real, "Attention for fine-grained categorization," *arXiv:1412.7054*, 2014.
- [21] S. Semeniuta and E. Barth, "Image classification with recurrent attention models," *Proc. 2016 IEEE Symposium Series on Computational Intelligence (SSCI)*, pp. 1-7, 2016.
- [22] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show attend and tell: Neural image caption generation with visual attention," *Proc. 32nd Int. Conf. Mach. Learn.*, pp. 2048-2057, 2015.
- [23] K. Xu, Y. Shi, L. Zheng, J. Zhang, M. Liu, H. Huang, H. Su, D. Cohen-Or and B. Chen, "3d attention-driven depth acquisition for object identification," *ACM Trans. Graphics*, vol. 35, no. 6, pp. 1-14, 2016.
- [24] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, "Learning representations by back-propagating errors," *Nature*, vol. 323, no. 6088, pp. 533-536, 1986.
- [25] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, no. 3-4, pp. 229-256, 1992.
- [26] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," *Proc. 2015 IEEE Conf. on Computer Vision and Pattern Recognition*, pp. 1912-1920, 2015.
- [27] D. Maturana and S. Scherer, "VoxNet: A 3D Convolutional Neural Network for real-time object recognition," *Proc. 2015 IEEE/RSJ International Conf. on Intelligent Robots and Systems (IROS)*, pp. 922-928, 2015.
- [28] A. Brock, T. Lim, J. Ritchie, and N. Weston, "Generative and Discriminative Voxel Modeling with Convolutional Neural Networks," *arXiv:1608.04236*, 2016.
- [29] Y. Li, S. Pirk, H. Su, C. R. Qi, and L. J. Guibas, "FPNN: Field Probing Neural Networks for 3D Data," *arXiv:1605.06240*, 2016.
- [30] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "PointNet: Deep Learning on Point Sets for 3D Classification and Segmentation," *arXiv:1612.00593*, 2016.
- [31] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "PointNet++: Deep Hierarchical Feature Learning on Point Sets in a Metric Space," *arXiv:1706.02413*, 2017.
- [32] P. S. Wang, Y. Liu, Y. X. Guo, C. Y. Sun, and X. Tong, "Octree-based Convolutional Neural Networks for 3D Shape Analysis," *ACM Trans. Graphics*, vol. 36, no. 4, pp. 1-11, 2017.

- [33] J. Li, B. Chen, and G. Lee, "SO-Net: Self-Organizing Network for Point Cloud Analysis," Proc. 2018 IEEE Conf. on Computer Vision and Pattern Recognition, pp. 9397-9406, 2018.
- [34] B. Shi, S. Bai, Z. Zhou, and X. Bai, "DeepPano: Deep panoramic representation for 3-D shape recognition," IEEE Signal Process. Lett., vol. 22, no. 12, pp. 2339-2343, 2015.
- [35] K. Sfikas, T. Theoharis and I. Pratikakis, "Exploiting the PANORAMA representation for convolutional neural network classification and retrieval," Proc. 2017 Eurographics Workshop on 3D Object Retrieval, 2017.
- [36] K. Sfikas, I. Pratikakis and T. Theoharis, "Ensemble of PANORAMA-based convolutional neural networks for 3D model classification and retrieval," Computers & Graphic, vol. 71, pp. 208-218, 2017.
- [37] 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, num. 7, pp. 1-11, 2015.
- [38] S. D. Roy, S. Chaudhury, and S. Banerjee, "Active recognition through next view planning: a survey," Pattern Recognition, vol. 37, no. 3, pp. 429-446, 2004.
- [39] K. Wu, R. Ranasinghe, and G. Dissanayake, "Active recognition and pose estimation of household objects in clutter," Proc. 2015 IEEE International Conference on Robotics and Automation (ICRA), pp. 4230-4237, 2015.
- [40] S. Yeung, O. Russakovsky, G. Mori, and L. Fei-Fei, "End-to-end learning of action detection from frame glimpses in videos," Proc. 2016 IEEE Conf. on Computer Vision and Pattern Recognition, pp. 2678-2687, 2016.
- [41] W. Zaremba and I. Sutskever, "Reinforcement learning neural Turing machines," arXiv:1505.00521, 2015.
- [42] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," Proc. 2016 IEEE Conf. on Computer Vision and Pattern Recognition, pp. 770-778, 2016.
- [43] S. Hochreiter and J. Schmidhuber, "Long short-term memory," Neural computation, vol. 9, no. 8, pp. 1735-1780, 1997.
- [44] D. Wierstra, A. Foerster, J. Peters, and J. Schmidhuber, "Solving deep memory POMDPs with recurrent policy gradients," Proc. International Conference on Artificial Neural Networks, pp. 697-706, 2007.
- [45] S. Hochreiter, "The vanishing gradient problem during learning recurrent neural nets and problem solutions," International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems, vol. 6, no. 2, pp. 107-116, 1998.
- [46] D. A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," arXiv:1511.07289, 2015.
- [47] T. Chen, M. Li, Y. Li, M. Lin, N. Wang, M. Wang, T. Xiao, B. Xu, C. Zhang, and Z. Zhang, "MXNet: A flexible and efficient machine learning library for heterogeneous distributed systems," arXiv:1512.01274, 2015.
- [48] V. Hegde and R. Zadeh, "Fusionnet: 3d object classification using multiple data representations," arXiv:1607.05695, 2016.
- [49] M. Jaderberg, K. Simonyan, and A. Zisserman. "Spatial transformer networks," Proc Advances in neural information processing systems, pp. 2017-2025, 2015.
- [50] W. Pei, T. Baltrušaitis, D. M. Tax, and L. P. Morency, "Temporal attention-gated model for robust sequence classification," Proc. 2017 IEEE Conf. on Computer Vision and Pattern Recognition, pp. 820-829, 2017.



Songle Chen received the Ph.D. degree from the Department of Computer Science and Technology of Nanjing University in 2015. Now he is a lecturer in Nanjing University of Posts and Telecommunications and a researcher of Jiangsu High Technology Research Key Laboratory for Wireless Sensor Networks. His research interests include Geometric Modeling, Realistic Rendering, Computer Animation, Deep Learning etc.



Lintao Zheng received his BS degree in applied mathematics from Xian Jiaotong University and MS degree in computer science from the National University of Defense Technology in 2013 and 2016, respectively. He is pursuing a doctorate in computer science at the National University of Defense Technology. His research interests mainly include computer graphics, deep learning, robot vision. He is a member of the ACM.



Yan Zhang received the Ph.D. degree from Jilin University in 2003 and finished his Post-doctoral researches in Nanjing University in 2006. She is now an associate professor of the Department of Computer Science and Technology in Nanjing University. Her research interests include Computer graphics, Image and Video Process and Computer Aided Design. Her research work can be found in her personal website: <http://cs.nju.edu.cn/zhangyan>.



Zhixin Sun is the professor and dean of School and Institute of Modern Posts, Nanjing University of Posts and Telecommunications. He received his Ph.D. degree in Nanjing University of Aeronautics and Astronautics, China in 1998 and worked as a post doctor in Seoul National University, South Korea between 2001 and 2002. His research area includes Computer Aided Design, computer vision, information security, computer networks, etc.



Kai Xu is an Associate Professor at the School of Computer, National University of Defense Technology, where he received his Ph.D. in 2011. He conducted visiting research at Simon Fraser University during 2008-2010, and Princeton University during 2017-2018. His research interests include geometry processing and geometric modeling, especially on data-driven approaches to the problems in those directions, as well as 3D-geometry-based computer vision. He has published over 60 research papers, including 21 SIGGRAPH/TOG papers. He organized two SIGGRAPH Asia courses and one Eurographics STAR tutorial. He is currently serving on the editorial board of Computer Graphics Forum, Computers & Graphics, and The Visual Computer. He also served as paper co-chair of CAD/Graphics 2017 and ICVRV 2017, as well as PC member for several prestigious conferences including SIGGRAPH Asia, SGP, PG, GMP, etc. His research work can be found in his personal website: www.kevinkaixu.net.