

GARNet: Global-Aware Multi-View 3D Reconstruction Network and the Cost-Performance Tradeoff

Zhenwei Zhu¹ Liying Yang¹ Xuxin Lin¹ Chaohao Jiang¹ Ning Li¹ Lin Yang² Yanyan Liang^{1,*}

¹Macau University of Science and Technology, Faculty of Innovation Engineering

²Kingsoft, Seasun Games

Abstract

Deep learning technology has made great progress in multi-view 3D reconstruction tasks. At present, most mainstream solutions establish the mapping between views and shape of an object by assembling the networks of 2D encoder and 3D decoder as the basic structure while they adopt different approaches to obtain aggregation of features from several views. Among them, the methods using attention-based fusion perform better and more stable than the others, however, they still have an obvious shortcoming — the strong independence of each view during predicting the weights for merging leads to a lack of adaption of the global state. In this paper, we propose a global-aware attention-based fusion approach that builds the correlation between each branch and the global to provide a comprehensive foundation for weights inference. In order to enhance the ability of the network, we introduce a novel loss function to supervise the shape overall and propose a dynamic two-stage training strategy that can effectively adapt to all reconstructors with attention-based fusion. Experiments on ShapeNet verify that our method outperforms existing SOTA methods while the amount of parameters is far less than the same type of algorithm, Pix2Vox++. Furthermore, we propose a view-reduction method based on maximizing diversity and discuss the cost-performance tradeoff of our model to achieve a better performance when facing heavy input amount and limited computational cost.

1. Introduction

3D reconstruction, a problem that involve the fields of computer vision and computer graphics, is considered as the core of many technologies such as computer-aided geometric design, computer animation, medical image processing, digital media and robotics. As a generation task, comparing to image restoration, lifting 2D images to 3D object is obviously an extremely difficult ill-posed inverse problem. According to the number of images as input, this task can be

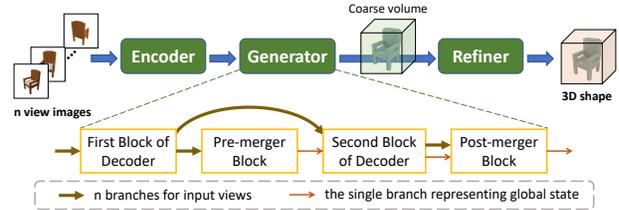


Figure 1. Our network consists of encoder, generator and refiner. The data flow in generator reflects the concept of global-aware.

divided into single-view reconstruction [15, 16, 18, 22, 25] and multi-view reconstruction [2–4, 6]. In this work, we focus on the deep-learning-based algorithms which reconstruct the shape of 3D objects with voxel representation from multiple images.

At present, most mainstream solutions assemble 2D encoder and 3D decoder as a basic framework and reshape the high-level features between them as a two-dimensional connection to establish the mapping between an image and a voxel. Nevertheless, there is still an important issue for multi-view reconstruction — how to aggregate the features from an arbitrary number of views.

In our investigation, there are four types of fusion strategies identified. [1, 32] adapt a pooling-based method to compress the feature map to a specific size using a pooling layer after concatenating the feature maps from all views. This dimensional collapse is too rough to avoid a massive loss of content. To make the fusion module learnable, 3D-R2N2 series [2, 9] utilize recurrent neural network (RNN)-based methods. The features from all views are regarded as a sequence and processed by a recurrent unit before the decoder. However, it indicates inconsistent predictions for different permutations. In addition, such methods are not suitable for too many views as input because of the limited long-term memory. To address these shortcomings, attention-based fusion approaches create a sub-network to predict the confidence score map of each view and merge features based on it. Both AttSets [3] and Pix2Vox series [4, 5] following this idea produce stable reconstructors. The former merges the features while the latter merges the

*Corresponding author

voxels restored from each view directly. Very recent, some researches [6, 7] use transformer structure for multi-view reconstruction. Utilizing the natural advantage, the fusion process is integrated into the encoder stage. They perform well for a large number of view inputs, but the restoration quality is terrible when there are few input images.

We consider that the attention-based fusion performs better and more stable than the others, however, it still has an obvious shortcoming. During predicting the score maps, the connection of branches only relies on a softmax layer without any learnable parameters, so that it cannot be adaptive to the global state and only trusts the memory of the network. To address this issue, we propose a global-aware multi-view 3D reconstruction network, named GAR-Net, which not only follows the attention-based fusion but also establishes the correlation between each branch and the global. Figure 1 illustrates the composition of the network and highlights the data flow in the generator, which reflects the concept of global-aware fusion.

In practical problems, we provide a large number of view images for better results, however, the model may not process due to the constraints of time consumption and computational complexity. To overcome this trouble, we propose a view-reduction method based on maximizing diversity and discuss the cost-performance tradeoff of our model to achieve a better performance when facing heavy input amount and limited computational cost.

In detail, the contributions are as follows:

- **Network architecture:** We propose a global-aware fusion approach that inherits the stability of the attention-based mergers and establishes the correlation between views to achieve better performance.
- **Loss function:** Precision and recall as quantitative indicators participate in the novel supervision to alleviate the weakness of expressing the difference between shapes brought by cross-entropy loss.
- **Training strategy:** We propose a dynamic two-stage training strategy that distinguishes the network processes for single-view and multi-view input and alternates training in these two situations randomly. This strategy can effectively adapt to all reconstructors using attention-based fusion.
- **Cost-performance tradeoff:** Facing heavy input amount, we provide a view-reduction method based on maximizing diversity to make the model achieve better performance under limited computational cost.

Furthermore, our models favorably against the SOTA methods [5, 6] in performance with fewer parameters than the same type of methods.

2. Related Works

- **Single-view 3D reconstruction.** In recent years, estimating 3D shape from a single view image is a hot topic. PointSetGeneration [15] generates 3D shape based on the representation of point clouds. Pixel2Mesh [16] represents the object by triangular mesh and process through a graph convolutional network (GCN). Voxel-based methods are rather widespread. [17] modifies voxel grids directly utilizing a 3DCNN. Using generative adversarial network (GAN) [20], 3DGAN [18] and 3DIWGAN [19] are proposed to solve the problem of 3D object generation, and the generator in these works can be converted to the single-view reconstructor by combining the variational auto-encoder (VAE) [21]. For the high-resolution results, OGN [22] adapts the octree representation to overcome the trouble of a huge memory budget and designs a network to process it directly, however, Matryoshka Networks [23] decompose a 3D shape into nested shape layers in a recursive manner. To bridge the gap between synthetic data and real-world data, DAREC [31] and VPAN [24] introduce the supervision on domain adaption during training. To supplement the missing information in the image, Mem3D [25] constructs a memory network to offer the priors information accumulated from the training set.
- **Multi-view 3D reconstruction.** SFM [26] and SLAM [27], the traditional reconstruction approaches, rely on matching features to establish the relationship across different views, but they have great restrictions on usage scenarios. Recently, deep-learning-based approaches are popular in multi-view 3D reconstruction, frequently without viewpoint labels. [28] leverages 2DCNN to predict dense point clouds representing the surface of 3D objects. In Pixel2Mesh++ [29], a coarse mesh can be improved iteratively with a series of deformations predicted by a GCN to form the final results. Representing by voxel, the methods focus on how to merge the features from several views. [1, 32] compress the concatenating features from all views using a maximum pooling layer. 3D-R2N2 series [2, 9] and LSM [30] receive the views one by one and extract the useful knowledge through a recurrent unit. EVoIT [6] and LegoFormer [7] exploit the advantage of transformer structure to realize the blending of information between various views in the encoder stage. As the most stable methods currently, AttSets [3] and Pix2Vox series [4, 5] apply the attention module on the multi-branch tasks, however, lack the exchange of information between branches.

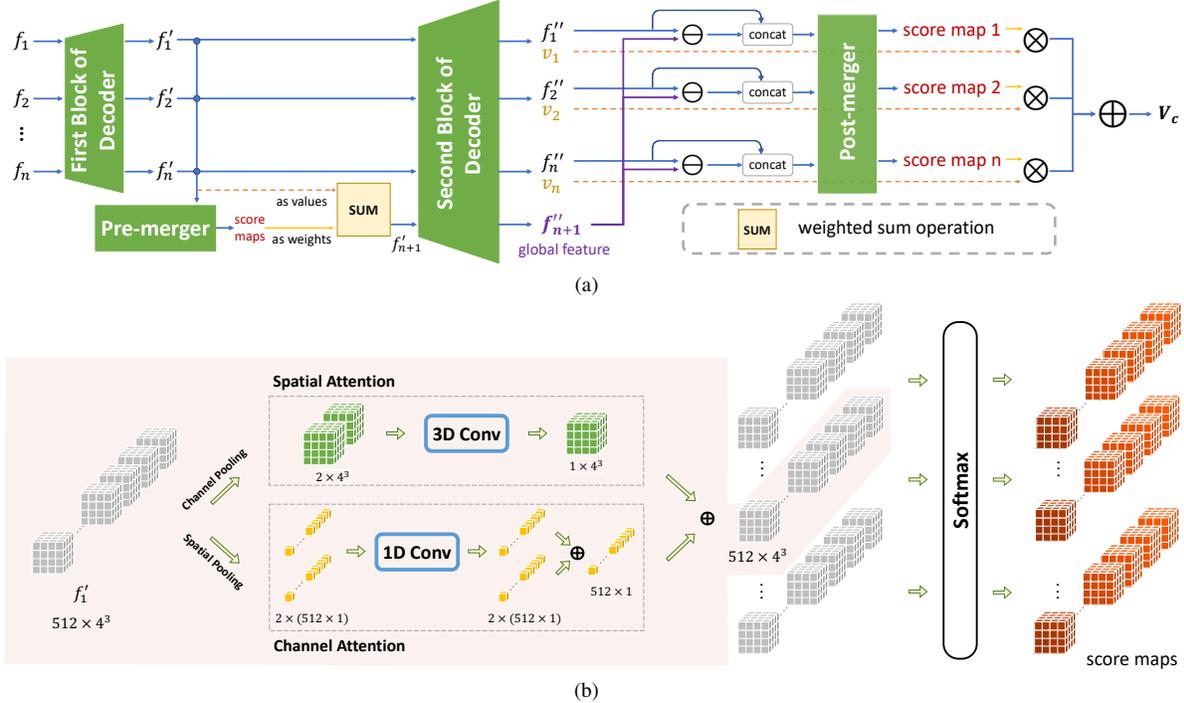


Figure 2. (a) is an overview of the generator that merges the features from different branches to reconstruct a coarse volume. It consists of two blocks of decoder, two fusion modules and some extra operations. Among them, the architecture of the pre-merger block, which predicts the score map for initial fusion, is shown as (b).

3. Methods

According to this task, the goal of our architecture is to bridge an uncertain number of RGB images $I = \{I_1, I_2, \dots, I_n\}$ with a size of $224 \times 224 \times 3$ to a binary voxel V representing the shape with a size of $32 \times 32 \times 32$. Our network consists of three parts: encoder, generator and refiner. To begin, the encoder extracts feature map f_i in parallel from each view image I_i . Then, the generator, which is formed by inserting two fusion blocks and a series of operations into a decoder, reconstructs a volume V_c from these features exploiting the global-aware attention-based fusion that discusses in section 3.1. However, this volume still has great potentials for making further progress. Inspired by [5], a refiner is adopted to modify V_c to the final output V . Mathematically, the complete network is defined as:

$$V = GARNet(I) = R(G(E(I_1, I_2, \dots, I_n))), \quad (1)$$

where E , G and R denote the encoder, generator and refiner respectively. Referring to prior experience, encoder and decoder used in generator share the same structure as [5]. The refiner is a novel network, named 3D-U-ResNet, that combine the advantages of ResNet [12] and U-Net [13]. It is a powerful structure with relatively lightweight and will be shown in the supplementary material.

3.1. Global-Aware Fusion

In this section, we elaborate on the processing in generator that implements the global-aware fusion. As shown in Figure 2a, the generator consists of four blocks, the two of which are split by the decoder and the others are pre-merger M_{pre} and post-merger M_{post} . The decoder includes four transposed convolutional layers with a kernel size of 4^3 and stride of 2 and one transposed convolutional layer with a kernel size of 1^3 and stride of 1. The first layer is regarded as the first block of decoder, referred to as D_1 , and the others compose the second block of decoder, referred to as D_2 .

The input of the generator is feature maps extracted from view images. Considering to utilize the spatial information of features, we prefer to process the initial fusion on 4D tensors. However, these features derived from encoder lack the spatial relationship. Therefore, the first block of decoder happening before the initial fusion is reasonable.

$$f'_i = D_1(E(I_i)) = D_1(f_i), \quad (2)$$

where f'_i is the output of D_1 on the branch of I_i .

Figure 2b shows the process of inferring the score map in pre-merger. Pre-merger employs two parallel attention mechanisms, a structure similar to the bottleneck attention module [10]. Both maximum pooling and average pool-

ing are utilized to compress features on channel and spatial to obtain a two-channel tensor and two vectors. The tensor is handled by a 3D convolution. For the high-dimensional vector, we replace the commonly fully connected layers with a 1D convolutional layer to maintain the structure lightweight inspired by [11]. After the parallel module, the branch feature map for initial fusion is created by extending and adding the channel perception feature map and the spatial perception feature map. The features of all views are combined and processed by a softmax layer to predict the corresponding normalized score maps for each branch and each grid in f' . Finally, the initial fusion feature map, referred to as f'_{n+1} , is the result of the addition of f' weighted by the score maps and becomes the beginning of the $(n + 1)$ th branch that need enter to D_2 like the other branches.

$$f'_{n+1} = \sum_i^n M_{pre}(f'_i) \cdot f'_i. \quad (3)$$

The second block of decoder will reconstruct a voxel v_i for each view. In addition, because the post-merger block needs to use the features of each branch to predict the score maps, we regard the concatenation of the feature maps from the last two layers of the decoder as the final features f''_i of the view image I_i .

$$(f''_i, v_i) = D_2(f'_i). \quad (4)$$

According to the features of the $(n + 1)$ branches, the post-merger block predicts the score maps of the first n branches corresponding to the restored volumes from the n view images. Since the actual meaning of the score map is the contribution rate of a branch to the entire, both the view and the global state are related to fusion logically. The input of post-merger for each branch is the concatenation of two feature maps: the feature map of the view f''_i and the deviation between f''_i and the global feature map f''_{n+1} from the $(n + 1)$ th branch. Due to only few channels of feature input, the post-merger, which consists of five 3D convolutional layers, has a more ordinary architecture than the pre-merger. Finally, following the same principle as pre-merger, the coarse volume is obtained after predicting the score maps with a softmax layer and fusing the reconstructed voxels from all view images weighted by the score maps.

$$V_C = \sum_i^n M_{post}(f''_i, f''_{n+1} - f''_i) \cdot v_i. \quad (5)$$

Summarizing the entire generator, we construct two attention-based fusion modules to implement global-aware. The first one merges all features to generate initial global features and then the second one makes full use of the global features to predict more reliable fusion weights for volumes.

3.2. Loss Function

The loss function is used to supervise both the coarse volume and the final volume. The binary cross-entropy (BCE) loss, which is commonly employed in previous works, is used to compare them to the ground truth respectively. Furthermore, we introduce precision and recall as new quantitative indicators to supervise the shape of the final volume. BCE only focuses on the classification of the cells but does not monitor overall shape. However, the precision and recall of occupied grids reflect the difference in shape. As a result, the combination can achieve the supervision to take both the local and the global into account. The complete loss function is defined as:

$$L = \underbrace{\alpha L_{BCE_{V_c}}}_{\text{for coarse volume}} + \underbrace{\beta L_{BCE_{V}} + \gamma L_{Recall} + \mu L_{Precision}}_{\text{for fine volume}}, \quad (6)$$

where $\alpha, \beta, \gamma, \mu$ indicate the weights of each part. $L_{BCE_{V_c}}$ and $L_{BCE_{V}}$ represent the BCE loss function of the two reconstruction results and the following two items are the recall loss function and the precision loss function defined as:

$$L_{Recall} = 1 - \frac{\sum_{i=1}^{32^3} p_i g t_i}{\sum_{i=1}^{32^3} g t_i}, \quad (7)$$

$$L_{Precision} = 1 - \frac{\sum_{i=1}^{32^3} p_i g t_i}{\sum_{i=1}^{32^3} p_i}, \quad (8)$$

where p and $g t$ denote the grids on the predicted result and ground truth.

4. Dynamic Two-Stage Training Strategy

At present, there are two strategies for training a reconstruction model using attention-based fusion. [3] proposes a feature-attention separate training (FASet) algorithm to split the training of the fusion module and the other parts into two stages. The first stage train all the parts except the fusion module using single-view input and the second stage only optimizes the fusion module using multi-view input. Pix2Vox series [4, 5] offer another two-stage training strategy with a similar first stage to FASet while the second stage fed with random numbers of input to train the complete network, in which the random number is updated for each epoch. However, when the random number is 1, the fusion modules cannot get an useful backpropagation gradient limited by the softmax layer, which is proved in [3]. So, the second stage in Pix2Vox training is not reasonable enough. In addition, we consider that training the model sufficiently using single image input firstly just like the two mentioned strategies may easily lead to overfitting for the single-view reconstruction task and limit the performance for multi-view.

To improve the performance of the attention-based fusion reconstructor, especially for multi-view input, we propose a general training strategy. First of all, the situations of single-view and multi-view as input are distinguished explicitly to perform different network processes. The propagation process does not go through the fusion module when the input is a single image, while for multi-view input, the parameters of the whole network participate. For each iteration, an integer less than or equal to the predefined value n_{max} will be randomly chosen as the number of view inputs. As a result, the two modes, single-view reconstruction and multi-view reconstruction, alternate training at a certain vague frequency. It not only utilizes the stability brought by the former to assist the train of the latter but also avoids model overfitting trending to the former task. The method to integrate the two training modes in this way is named dynamic two-stage training strategy and its effectiveness is verified in Section 6.2.2.

5. Cost-Performance Tradeoff via View-Reduction

We have established a robust and stable multi-view 3D reconstruction network so far. In practice, users can employ more view images as input to achieve better reconstruction results. However, increasing input images means greater computational complexity and more time consumption. In some cases, the acceptable time for inference is strictly limited. It is necessary to provide a cost-performance tradeoff analysis for generating relatively high-quality results with limited computational cost. Considering the multi-view reconstruction accuracy relies on the diverse viewpoints information, we propose a view-reduction approach based on maximizing diversity to remove some branches and control the computational cost while preserving performance. Intuitively, a combination of views with a considerable disparity in viewpoints taking place of similar views can provide more diverse information for restoring a reliable volume. As a result, we attempt to find a combination of images with roughly complementary information from all view input to reduce the branches. The approach is based on isometric mapping training and maximizing diversity selection.

5.1. Isometric Mapping

To extract the diverse viewpoint information via reducing the redundant features, we assume a low-dimensional manifold space in which each point represents observation position information matching to a view image of the object, and the Euclidean distance between two points is consistent with the difference of the viewpoints. Derived from the space, it is convenient to pick a subset of points with a widely scattered distribution that corresponds to an image combination with relatively complementary information.

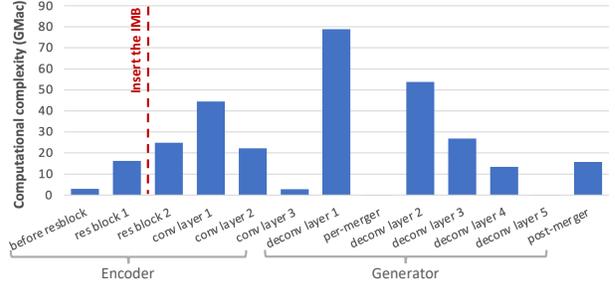


Figure 3. The distribution of computational complexity in encoder and generator.

Therefore, we build an isometric mapping block (IMB) with a simple structure composed of two paralleled pooling layers and three MLPs (multilayer perceptron) for dimensionality reduction of encoded features, and then insert it into the reconstruction model to obtain a more compact representation preserving information of viewpoints. Such that we can use fewer view branches to reconstruct the volume based on maximizing diversity of viewpoints information. Figure 3 records the distribution of computational complexity before the refiner when facing a heavy input amount (24 views). We expect to assign fewer operations before IMB since the quantity of calculation can only be reduced after that by eliminating branches. As a result, it is reasonable to use the output of the first residual block in the encoder as the input of IMB. Training IMB still adopts the data without location and direction information of views. As mentioned by [4], the score maps for volumes predicted by the fusion module can be thought of as a representation of the visible parts from a viewpoint learned by the network adaptively. The weights allocated to the visible parts will presumably be higher. Consequently, the difference between the score maps predicted by the post-merger block in our model can be used to distinguish the difference between the visible parts, i.e. the relationship of viewpoints.

To train IMB, a manifold learning algorithm similar to [14] is established for isometric mapping the low-dimensional points in Euclidean space to the score maps for different views of one object in L1 space. The L1 loss utilized to supervise is defined as:

$$L_{IMB} = \left\| \sum_i^n \sum_j^n (\|P_i - P_j\|_{L^2} - \|Q_i - Q_j\|_{L^1}) \right\|_{L^1}, \quad (9)$$

where P represents the points in low-dimensional manifold space corresponding to view images of an object and Q denotes their matching score maps.

	1 view	2 views	3 views	4 views	5 views	8 views	12 views	16 views	20 views
3D-R2N2 [2]	0.560 / 0.351	0.603 / 0.368	0.617 / 0.372	0.625 / 0.378	0.634 / 0.382	0.635 / 0.383	0.636 / 0.382	0.636 / 0.382	0.636 / 0.383
AttSets [3]	0.642 / 0.395	0.662 / 0.418	0.670 / 0.426	0.675 / 0.430	0.677 / 0.432	0.685 / 0.444	0.688 / 0.445	0.692 / 0.447	0.693 / 0.448
Pix2Vox++ [5]	0.670 / 0.436	0.695 / 0.452	0.704 / 0.455	0.708 / 0.457	0.711 / 0.458	0.715 / 0.459	0.717 / 0.460	0.718 / 0.461	0.719 / 0.462
EvoIT [6]	- / -	- / -	- / -	0.609 / 0.358	- / -	0.698 / 0.448	0.720 / 0.475	0.729 / 0.486	0.735 / 0.492
Legoformer [7]	0.519 / 0.282	0.644 / 0.392	0.679 / 0.428	0.694 / 0.444	0.703 / 0.453	0.713 / 0.464	0.717 / 0.470	0.719 / 0.472	0.721 / 0.472
GARNet	0.673 / 0.418	0.705 / 0.455	0.716 / 0.468	0.722 / 0.475	0.726 / 0.479	0.731 / 0.486	0.734 / 0.489	0.736 / 0.491	0.737 / 0.492
GARNet+	0.655 / 0.399	0.696 / 0.446	0.712 / 0.465	0.719 / 0.475	0.725 / 0.481	0.733 / 0.491	0.737 / 0.498	0.740 / 0.501	0.742 / 0.504

Table 1. Evaluation and comparison of the performance on ShapeNet using IoU / F-Score@1%. The best results are highlighted in bold.

	1 view	2 views	3 views	4 views	5 views
Setup 1	0.670	0.695	0.704	0.708	0.711
Setup 2	0.6693	0.6990	0.7090	0.7138	0.7172
Setup 3	0.6707	0.7026	0.7136	0.7187	0.7223
Setup 4	0.6697	0.7023	0.7137	0.7195	0.7235
Setup 5	0.6725	0.7047	0.7160	0.7217	0.7255
Setup 6	0.6551	0.6958	0.7117	0.7193	0.7246
	8 views	12 views	16 views	20 views	
Setup 1	0.715	0.717	0.718	0.719	
Setup 2	0.7217	0.7245	0.7262	0.7269	
Setup 3	0.7280	0.7308	0.7324	0.7332	
Setup 4	0.7289	0.7319	0.7336	0.7345	
Setup 5	0.7312	0.7340	0.7357	0.7368	
Setup 6	0.7331	0.7373	0.7400	0.7415	

Table 2. The ablation experiments on ShapeNet about dynamic two-stage training strategy, 3D-U-ResNet as the refiner network, global-aware fusion, precision-recall loss function and 8-view input.

5.2. Diversity Maximization

The problem is now formulated as reducing N branches of views to n . We obtain N low-dimensional points via IMB. Using farthest point sampling (FPS), n points can be selected and only their corresponding branches are retained to continue to finish the rest of the reconstruction network.

FPS is a sampling method with uniform and wide coverage so that the selected points maximize the diversity of the entire point set in a limited capacity. These points correspond to a view combination with a great variation in viewpoints. Meanwhile, the score maps, which indicate the contribution of views to different voxel grids, are directly tied to the meaning of these points. Either retaining branches with widely varying viewpoints or larger differences of contribution weights distribution imply maximizing diversity.

Thus, the view-reduction method for cost-performance tradeoff by retaining the view combination with more diverse information for reconstruction is realized and the effectiveness will be verified in Section 6.3.

6. Experiments

- **Dataset.** We evaluate our reconstruction network on the ShapeNet [8] dataset using both Intersection of Union (IoU) and F-Score@1% [5, 33] as the metric.

Following [2], only a subset of ShapeNet including 13 categories and 43,783 3D objects with 24 randomly view images for each are used in our experiments.

- **Implementation Details.** We adopt an Adam optimizer [34] with $\beta_1 = 0.9$ and $\beta_2 = 0.999$ to train our multi-view reconstruction network with a batch size of 32 for 200 epochs, with 140 epochs using only BCE loss function and then 60 epochs using the complete loss function as mentioned previously. The weights in Equation 6 are set as $\alpha = \beta = 10$ and $\gamma = \theta = 0.5$. The learning rate is $1e-3$ initially and reduce to half after [40, 60, 80, 100, 140, 180] epochs sequentially. For the property of dynamic two-stage training strategy, the parameters in the fusion module are optimized less frequently than the others. Thus, setting a slightly higher learning rate for the fusion module part can archive a better performance. Eventually, we provide two models respectively setting the maximum number of input views to 3 and 8 during training, named GARNet and GARNet+. The fixed threshold for binarizing the probabilities is set as 0.3. It takes about 2 days to train GARNet on 1 Tesla V100 and about 3 days to train GARNet+ on 2 Tesla V100. The IMB for view-reduction is extra trained for 15 epochs relied on a converged and frozen reconstruction network, and uses 40 objects with 24 views for each iteration. The learning rate is set to $2e-3$ and decreased by 0.1 every 5 epochs.

6.1. Multi-View Reconstruction Results

The performance of our models for multi-view reconstruction is evaluated and compared with existing SOTA methods. As shown in Table 1, our proposed GARNet already dominates in almost all metrics. Not only does it performs consistently well when few views as input just like the other reconstructors using attention-based fusion, but also it breaks the advantage of the transformer method EvoIT [6] when a large number of views as input. Furthermore, we also provide GARNet+, which performs better in multi-view reconstruction tasks while sacrificing acceptable performance for single-view. The two models are identical in structure and both of them outperform existing SOTA methods while only using about 69% of parameters of the same type of algorithm, Pix2Vox++ [5].

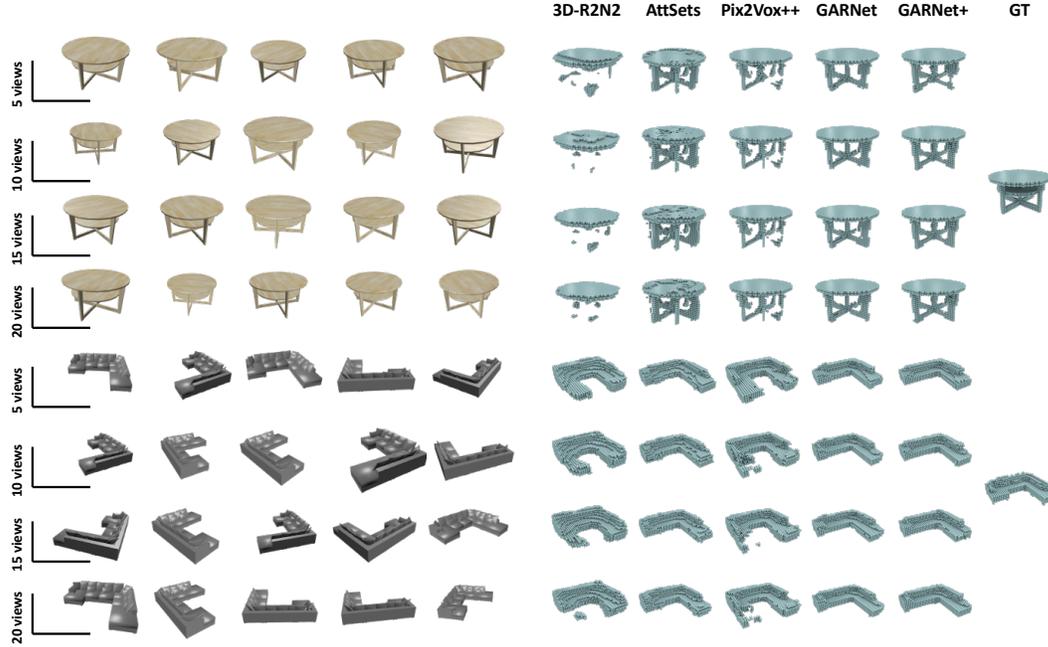


Figure 4. Multi-view reconstruction results on the test set of ShapeNet when facing 5 views, 10 views and 15 views as input.

CA layer	1 view	2 views	3 views	4 views	5 views
FC	0.6652	0.6991	0.7111	0.7167	0.7203
1D Conv	0.6697	0.7023	0.7137	0.7195	0.7235

CA layer	8 views	12 views	16 views	20 views
FC	0.7259	0.7295	0.7313	0.7322
1D Conv	0.7289	0.7319	0.7336	0.7345

Table 3. Comparison of performance evaluated by IoU when using fully connection (FC) layers and a 1D convolutional layer in channel attention (CA) module of pre-merger. Experiments on ShapeNet using only BCE loss function and setting the upper limit of input views to 3 during training.

As examples, Figure 4 shows several reconstruction results. Comparing to the other methods, for table restoration, our models present a smoother plane and depict the links between the legs more accurately. In addition, our results for the sofa are also stable and reasonable, however, Pix2Vox++ produces a part of incorrect and confusing voxel girds. Without the global-aware, mistakes in a branch will seriously affect the final result. We exploit the fusion at the feature level, which is relatively insensitive to local perception, as the global information to bring the network a certain ability for self-correction.

6.2. Ablation Experiments

We use Pix2Vox++ [5] as the baseline. To design ablation experiments, our proposed approaches, which include dynamic two-stage training strategy, 3D-U-ResNet as the refiner network, global-aware fusion, and precision-recall

loss function, are applied one by one based on it to convert the model to GARNet. Table 2 shows the results. Specifically, the setups of these experiments are as follows:

- **Setup 1:** Baseline (Pix2Vox++)
- **Setup 2:** Baseline + Dynamic two-stage training
- **Setup 3:** Baseline + Dynamic two-stage training+3D-U-Resnet
- **Setup 4:** Baseline + Dynamic two-stage training + 3D-U-Resnet + Global-aware fusion
- **Setup 5:** Baseline + Dynamic two-stage training + 3D-U-Resnet + Global-aware fusion + P-R loss (GARNet)
- **Setup 6:** Baseline + Dynamic two-stage training + 3D-U-Resnet + Global-aware fusion + P-R loss + 8 views (GARNet+)

It verifies that these methods have a positive impact on the reconstruction algorithm. In addition, GARNet+ training with 8 views input performs well when facing a large number of views input. In this section, we will discuss the relevant verification works about pre-merger structure and training strategy in detail.

6.2.1 Pre-Merger Block

We replace the fully connection layers with a 1D convolutional layer in the channel attention module of pre-merger

Training Strategy	1 view	2 views	3 views	4 views	5 views
FASet [3]	0.6733	0.6994	0.7082	0.7126	0.7154
Pix2Vox [4]	0.6708	0.7019	0.7122	0.7173	0.7203
Ours	0.6697	0.7023	0.7137	0.7195	0.7235
Training Strategy	8 views	12 views	16 views	20 views	
FASet	0.7195	0.7215	0.7228	0.7232	
Pix2Vox	0.7250	0.7276	0.7290	0.7297	
Ours	0.7289	0.7319	0.7336	0.7345	

Table 4. Experiments on our architecture using only BCE loss function and setting the upper limit of input views to 3 during training to compare our proposed dynamic two-stage training strategy with the two existing strategies on ShapeNet using IoU.

block. Heavy parameters are not expected for a branch of weights calculation, because their lack of direct supervision will make training more difficult. In addition, channels of a vector are not strongly related to each other in the channel attention module. As a result, a simpler and lighter structure bring better performance, as shown in Table 3.

The pre-merger block plays a critical role with only hundreds of parameters. It improves the performance of the model effectively by allowing data to flow across branches. To establish this connection is the most important responsibility of it, which increases countless potential network branches to promote the learning ability of the network for merging. For the feature representation, the decoder with a large number of parameters also assists it to generate a better global representation, so the network never lacks parameters but connections.

6.2.2 Training Strategy

As aforementioned, the dynamic two-stage training strategy is more reasonable than the previous methods for the models using attention-based fusion. As shown in Table 4, we employ these three strategies to train our network. As a result, our method has a distinct advantage in multi-view reconstruction tasks. For single-view reconstruction, our performance is slightly lower than theirs, but it is roughly the same. However, their models are overfitting for single-view reconstruction since training with single-view inputs sufficiently in the first stage, which is detrimental to the generalization ability of the model.

In addition, we explore the influence of the maximum number of input views during training for our strategy. Comparing GARNet and GARNet+, setting a higher upper limit with more memory consumption results in better performance on multi-view reconstruction, while losing the performance of single-view reconstruction due to the lower frequency of single-view input during training.

6.3. Cost-Performance Tradeoff

First of all, it is necessary to verify that our view-reduction method can retain the view combination includ-

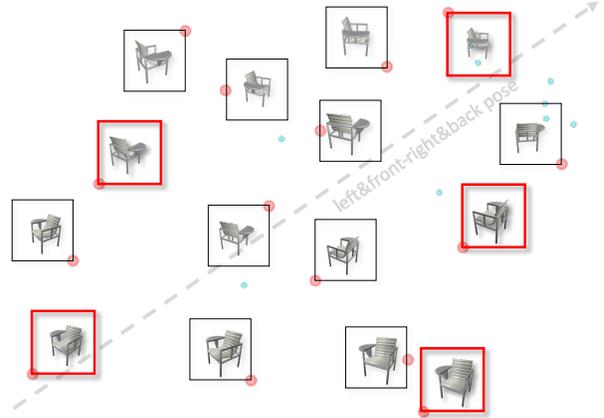


Figure 5. The distribution of the points in the 2D manifold space mapped from the view images of an object by a trained IMB.

	View-Reduction	24 → 3	24 → 4	24 → 5
GARNet	Random	0.7163	0.7221	0.7254
	FPS directly	0.7166	0.7223	0.7256
	PCA + FPS	0.7180	0.7243	0.7279
	IMB + FPS (ours)	0.7225	0.7269	0.7295
GARNet+	Random	0.7120	0.7201	0.7246
	FPS directly	0.7125	0.7205	0.7258
	PCA + FPS	0.7138	0.7223	0.7281
	IMB + FPS (ours)	0.7200	0.7263	0.7305

Table 5. Comparison of the reconstruction performance based on the specified quantity of images reduced from 24 views by random sampling, FPS directly, PCA + FPS and our method, which is evaluated on test set of ShapeNet using IoU.

ing more diverse information. The different views of an object are mapped to a 2D manifold space by IMB and shown in Figure 5. There is a certain correlation between the observation position of the chair and the location of the corresponding points in the space. It means that IMB extracts the features about viewpoints from the images. In addition, the 5 objects marked with red frames are selected by FPS and their viewpoints are obviously different from each other. Thus, view-reduction based on isometric mapping and diversity maximization is reliable. Note that, the displayed distribution is not ideally perfect, because the 2D vector cannot adequately capture the difference across viewpoints. It is merely for visualization clarity and we set a higher dimension (empirically set 5) for experiments.

In the experiments, we train IMB for GARNet and GARNet+ respectively. Table 5 illustrates that using our method can produce better results comparing to random sampling, using FPS directly or using FPS after PCA when reconstructing a volume using a specified number of images selected from 24 views. According to the statistics, using our view-reduction approach to retain 5 branches can achieve about 98% precision of the results created based on all 24 views while saving about 70% of computational cost.

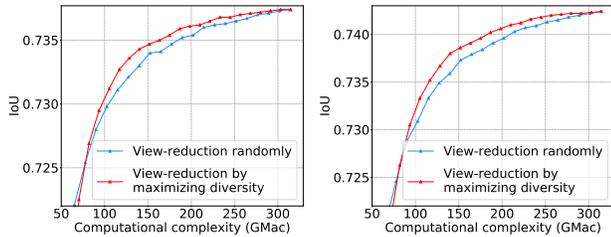


Figure 6. The cost-performance tradeoff that uses the maximizing diversity method and random sampling for view-reduction. The left one is the curve for GARNet and the right is for GARNet+.

For the cost-performance tradeoff, we have to analyze the relationship between specific computational cost and performance. When achieving the maximizing diversity selection, each view input must be processed by the network before the IMB. It is not required for view-reduction randomly. As a result, our method is not suitable for situations where the complexity is limited to an extremely low level. However, our view-reduction method still has applicable scenarios. Figure 6 presents the comparison of maximizing diversity and random sampling for view-reduction on GARNet and GARNet+. When the computational complexity is greater than 80 GMac (multiply and accumulate) is allowed, our method takes advantage. In practical problems, such a chart can help us determine a better course of action. Furthermore, the size of the reconstructed volume is only 32^3 with a low resolution because of the limitation of the dataset. If the model includes more decoder layers for a higher resolution result, view-reduction by maximizing diversity will obviously lead to a greater benefit.

7. Conclusion

In this paper, we propose a multi-view 3D reconstruction network using global-aware fusion, an advanced attention-based method, to favorably against the SOTA methods in performance. Furthermore, the cost-performance tradeoff is discussed for facing practical problems. In future work, we expect that these methods designed for the multi-branch network can be used to solve the multi-view reconstruction problem based on the other 3D representations and achieve high-resolution results.

References

- [1] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, “Multi-view convolutional neural networks for 3d shape recognition,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 945–953.
- [2] C. B. Choy, D. Xu, J. Gwak, K. Chen, and S. Savarese, “3d-r2n2: A unified approach for single and multi-view 3d object reconstruction,” in *European conference on computer vision*. Springer, 2016, pp. 628–644.
- [3] B. Yang, S. Wang, A. Markham, and N. Trigoni, “Robust attentional aggregation of deep feature sets for multi-view 3d reconstruction,” *International Journal of Computer Vision*, vol. 128, no. 1, pp. 53–73, 2020.
- [4] H. Xie, H. Yao, X. Sun, S. Zhou, and S. Zhang, “Pix2vox: Context-aware 3d reconstruction from single and multi-view images,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2690–2698.
- [5] H. Xie, H. Yao, S. Zhang, S. Zhou, and W. Sun, “Pix2vox++: multi-scale context-aware 3d object reconstruction from single and multiple images,” *International Journal of Computer Vision*, vol. 128, no. 12, pp. 2919–2935, 2020.
- [6] D. Wang, X. Cui, X. Chen, Z. Zou, T. Shi, S. Salcudean, Z. J. Wang, and R. Ward, “Multi-view 3d reconstruction with transformer,” *arXiv preprint arXiv:2103.12957*, 2021.
- [7] F. Yagubbayli, A. Tonioni, and F. Tombari, “Legoformer: Transformers for block-by-block multi-view 3d reconstruction,” *arXiv preprint arXiv:2106.12102*, 2021.
- [8] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, “3d shapenets: A deep representation for volumetric shapes,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1912–1920.
- [9] T. Ma, P. Kuang, and W. Tian, “An improved recurrent neural networks for 3d object reconstruction,” *Applied Intelligence*, vol. 50, no. 3, pp. 905–923, 2020.
- [10] J. Park, S. Woo, J.-Y. Lee, and I. S. Kweon, “Bam: Bottleneck attention module,” *arXiv preprint arXiv:1807.06514*, 2018.
- [11] Q. Wang, B. Wu, P. Zhu, P. Li, and Q. Hu, “Eca-net: Efficient channel attention for deep convolutional neural networks,” in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [12] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [13] O. Ronneberger, P. Fischer, and T. Brox, “U-net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [14] J. B. Tenenbaum, V. De Silva, and J. C. Langford, “A global geometric framework for nonlinear dimensionality reduction,” *science*, vol. 290, no. 5500, pp. 2319–2323, 2000.
- [15] H. Fan, H. Su, and L. J. Guibas, “A point set generation network for 3d object reconstruction from a single image,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 605–613.
- [16] N. Wang, Y. Zhang, Z. Li, Y. Fu, W. Liu, and Y.-G. Jiang, “Pixel2mesh: Generating 3d mesh models from single rgb images,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 52–67.

- [17] R. Girdhar, D. F. Fouhey, M. Rodriguez, and A. Gupta, “Learning a predictable and generative vector representation for objects,” in *European Conference on Computer Vision*. Springer, 2016, pp. 484–499.
- [18] J. Wu, C. Zhang, T. Xue, W. T. Freeman, and J. B. Tenenbaum, “Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling,” in *Proceedings of the 30th International Conference on Neural Information Processing Systems*, 2016, pp. 82–90.
- [19] E. J. Smith and D. Meger, “Improved adversarial systems for 3d object generation and reconstruction,” in *Conference on Robot Learning*. PMLR, 2017, pp. 87–96.
- [20] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” *Advances in neural information processing systems*, vol. 27, 2014.
- [21] D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” in *ICLR*, 2014.
- [22] M. Tatarchenko, A. Dosovitskiy, and T. Brox, “Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2088–2096.
- [23] S. R. Richter and S. Roth, “Matryoshka networks: Predicting 3d geometry via nested shape layers,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 1936–1944.
- [24] Q. Feng, Y. Luo, K. Luo, and Y. Yang, “Look, cast and mold: Learning 3d shape manifold from single-view synthetic data,” *arXiv preprint arXiv:2103.04789*, 2021.
- [25] S. Yang, M. Xu, H. Xie, S. Perry, and J. Xia, “Single-view 3d object reconstruction from shape priors in memory,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 3152–3161.
- [26] O. Özyeşil, V. Voroninski, R. Basri, and A. Singer, “A survey of structure from motion*,” *Acta Numerica*, vol. 26, pp. 305–364, 2017.
- [27] J. Fuentes-Pacheco, J. Ruiz-Ascencio, and J. M. Rendón-Mancha, “Visual simultaneous localization and mapping: a survey,” *Artificial intelligence review*, vol. 43, no. 1, pp. 55–81, 2015.
- [28] C.-H. Lin, C. Kong, and S. Lucey, “Learning efficient point cloud generation for dense 3d object reconstruction,” in *proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [29] C. Wen, Y. Zhang, Z. Li, and Y. Fu, “Pixel2mesh++: Multi-view 3d mesh generation via deformation,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 1042–1051.
- [30] A. Kar, C. Häne, and J. Malik, “Learning a multi-view stereo machine,” *arXiv preprint arXiv:1708.05375*, 2017.
- [31] P. O. Pinheiro, N. Rostamzadeh, and S. Ahn, “Domain-adaptive single-view 3d reconstruction,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7638–7647.
- [32] D. Paschalidou, O. Ulusoy, C. Schmitt, L. Van Gool, and A. Geiger, “Raynet: Learning volumetric 3d reconstruction with ray potentials,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3897–3906.
- [33] M. Tatarchenko, S. R. Richter, R. Ranftl, Z. Li, V. Koltun, and T. Brox, “What do single-view 3d reconstruction networks learn?” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3405–3414.
- [34] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.