

NormalNet: Learning-based Normal Filtering for Mesh Denoising

Wenbo Zhao, Xianming Liu, *Member, IEEE*, Yongsen Zhao, Xiaopeng Fan, *Senior Member, IEEE* Debin Zhao *Member, IEEE*,

Abstract—Mesh denoising is a critical technology in geometry processing that aims to recover high-fidelity 3D mesh models of objects from their noise-corrupted versions. In this work, we propose a learning-based normal filtering scheme for mesh denoising called *NormalNet*, which maps the guided normal filtering (GNF) into a deep network. The scheme follows the iterative framework of filtering-based mesh denoising. During each iteration, first, the voxelization strategy is applied on each face in a mesh to transform the irregular local structure into the regular volumetric representation, therefore, both the structure and face normal information are preserved and the convolution operations in CNN(Convolutional Neural Network) can be easily performed. Second, instead of the guidance normal generation and the guided filtering in GNF, a deep CNN is designed, which takes the volumetric representation as input, and outputs the learned filtered normals. At last, the vertex positions are updated according to the filtered normals. Specifically, the iterative training framework is proposed, in which the generation of training data and the network training are alternately performed, whereas the ground truth normals are taken as the guidance normals in GNF to get the target normals. Compared to state-of-the-art works, *NormalNet* can effectively remove noise while preserving the original features and avoiding pseudo-features.

Index Terms—Mesh denoising, convolutional neural networks, normal filtering, guided normal filtering, voxelization



1 INTRODUCTION

RECENTLY, a demand for high-fidelity 3D mesh models of real objects has appeared in many domains, such as computer graphics, geometric modelling, computer-aided design and the movie industry. However, due to the accuracy limitations of scanning devices, raw mesh models are inevitably contaminated by noise, leading to corrupted features that profoundly affect the subsequent applications of meshes. Hence, mesh denoising has become an active research topic in the area of geometry processing.

Mesh denoising is an ill-posed inverse problem. The nature of mesh denoising is to smooth a noisy surface, while concurrently preserving the real object features without introducing unnatural geometric distortions. Mesh denoising is a challenging task, especially for cases with large and dense meshes and with high noise levels. The key to the success of mesh denoising is to differentiate the actual geometry features, such as the localized curvature changes and small-scale details, and the noise generated by scanners. The literature contains rich work on mesh denoising, including filtering-based [1], [2], [3], [4], [5], [6], [7], feature-extraction-based [8], [9], optimization-based [10], [11], and similarity-based [12], [13], [14]. Among these methods, the guided-normal based scheme has become popular in recent years [3], [4], [5], [6], [7], which follows the iterative framework of filtering-based mesh denoising. During each iteration, the guidance normals are derived and used in filtering first. Then the vertex positions are updated according to the filtered normals. This approach performs mesh denoising either by building guidance normals with manually designed methods [3], [6], [7], or by introduc-

ing additional information to improve the performance of guidance normals [4], [5]. The schemes in [3], [4], [5] perform well on synthetic meshes with simple structures. However, the methods of generating guidance normals in [3], [4], [5] are based on finding consistent patches with fixed shapes, therefore cannot handle complex structures well such as narrow edges and corners. To overcome this problem, Li *et al.* [6] propose to generate the guidance normals by the corner-aware and edge-aware neighbourhood. Zhao *et al.* [7] employ the graph-cut to generate piece-wise smooth patches and build guidance normals on them. These schemes perform well on synthetic meshes with complex features. However, the main idea of these schemes is finding consistent patches according to the face normal difference, and the structure information has not been fully utilized. For scanned meshes, which contain manifold kinds of noise and more complex shapes, such as serrated noise (Fig. 11), stair-stepping noise (Fig. 14) and irregular edges (Fig. 10), the face normal difference of noisy faces is so large that it is difficult for [6], [7] to distinguish noise and features, resulting in either introducing pseudo-features or over-smooth. As shown in the experimental comparisons, even the state-of-the-art schemes [3], [7] cannot handle these cases well.

In the counterpart 2D image denoising, deep-learning-based strategies, such as [15], [16], [17], have been widely applied and achieved great success. However, with respect to mesh denoising, to the best of our knowledge, no studies follow this line of research. One main reason preventing the usage of convolutional neural network(CNN) in mesh denoising is that, in contrast to the regular grid structure of 2D images, meshes have irregular structures. Therefore, it is not straightforward to apply the regular 3D convolutional operations in CNNs to a mesh. Another reason may be the difficulty of selecting an efficient denoising strategy for

W. Zhao, X. Liu, Y. Zhao, X. Fan and D. Zhao are with the School of Computer Science and Technology, Harbin Institute of Technology, Harbin, 150001, China. e-mail: {wzbzhao, csxm, 18S103199, fxp, dbzhao}@hit.edu.cn.

CNN to mimic.

In this work, we propose a learning-based normal filtering scheme for mesh denoising called *NormalNet*, which maps the guided normal filtering (GNF) [3] into a deep network. In particular, *NormalNet* follows the iterative framework of GNF, as shown in Fig. 1. During each iteration, to overcome the difficulty in using CNNs on meshes and exploit both the structure and face normal information, first, the voxelization strategy is applied on each face in a mesh to convert the irregular local structure into the regular volumetric representation. Second, a deep CNN is designed, which takes the volumetric representation as input, and outputs the learned filtered normals. All CNNs share the same workflow: three residual blocks, one max-pooling layer and four fully connected layers, of which the fourth layer outputs the filtered normals. At last, the vertex positions are updated according to the filtered normals. Moreover, we propose an iterative training framework for *NormalNet*, in which the generation of training data and the training of CNNs are performed alternately, whereas the ground truth normals are taken as the guidance normals in GNF to get the target normals. Compared to the state-of-the-art schemes, *NormalNet* can effectively generate more accurate filtering results and remove noise while preserving the original features and avoiding pseudo-features.

The rest of this paper is organized as follows. In the following section, we briefly summarize the related work. The proposed *NormalNet* is introduced in Section 3. In Section 4, the training of *NormalNet* is elaborated. Experimental results are presented in Section 5. Finally, Section 6 concludes the paper.

2 RELATED WORK

In this section, we briefly review the related work on the filtering-based mesh denoising and the neural-network-based 3D model processing.

2.1 Filtering-based mesh denoising

Owing to the edge-preserving property of the bilateral filter, researchers have made many attempts to adopt the bilateral filtering in mesh denoising [18], [19], [20]. Nevertheless, the photometric weights in the bilateral filter cannot be estimated accurately from a noise-corrupted mesh. The joint bilateral filter [21], in which the photometric weights are computed from a reliable guidance image, is proposed to improve the capability of bilateral filtering. Inspired by this idea, Zhang *et al.* [3] propose the guided normal filtering, in which the guidance information is obtained as the average normal of a local patch. This scheme works well with respect to feature preservation but cannot achieve satisfactory results in regions with complex shapes and sometimes introduces pseudo-features. To overcome the problems of [3], in the subsequent work [6], the guidance normals are computed by the corner-aware neighbourhood, which is adaptive to the shapes of corners and edges.

Recently, there have been increasing efforts to exploit the geometric attributes for mesh denoising. In [22], the normal filtering is performed by means of a total variation, which assumes the normal change is piecewise constant. Wei *et*

al. [2] propose to cluster faces into piecewise-smooth patches and refine face normals with the help of vertex normal fields. In [10], a differential edge operator is proposed and the L0 minimization is employed to remove noise while preserving the sharp features. Further more, Lu *et al.* [23] apply an additional vertex filtering before the L1-median face normal filtering, which proves to be capable of handling high noise levels and noise distributed in a random direction. However, feature information, such as edges and corners with less noise, may be blurred due to prefiltering. In [24], the Tukey bi-weight similarity function is proposed to replace the similarity function in the computation of bilateral weights; in addition, an edge-weighted Laplace operator is introduced for vertex updating to reduce face normal flips. In [7], the graph-based feature detection is employed to construct accurate guidance normals; however, this method may introduce pseudo-features when the shape of the noise is complex, which is common in scanned models.

2.2 Neural-Network-based 3D model processing

Driven by the great success of deep learning in image processing, researchers in graphics are also attempting to employ deep neural networks for 3D model processing. However, due to the property of irregular connectivity, processing 3D models with neural networks remains challenging. Numerous works have focused on transforming 3D models into regular data. For instance, in [25], [26], 3D models are represented by 2D rendered images and panoramic views. Furthermore, some studies [27], [28], [29], [30] have employed voxelization to transform models into regular 3D data. Moreover, in [31], [32], [33], [34], meshes are represented in the spectral or spatial domain for further processing.

In addition to these transform-based techniques, the direct application of neural networks to irregular data has also been extensively studied for point cloud data. PointNet [35] is one of the first network architectures that can handle point cloud data. Subsequently, PointNet++ [36] and the dynamic graph CNN [37] are proposed to improve the network capability. Some attempts have been made to organize point clouds into structures. In [38], a kd-tree is constructed on a point cloud and is further used as the input of a neural network. A similar idea is presented in [39], where the points are organized by an octree. Additional works [40], [41], [42] focus on surface reconstruction, denoising and removing outliers.

Notably, in [43], Wang *et al.* propose the filtered facet normal descriptor and model it with neural networks, however, these networks are not convolutional and only take face normal information into considered. In [44], the edge-based convolution and pooling operations are defined which can be directly on the constructs of the mesh.

3 THE FRAMEWORK OF NORMALNET

In this section, we introduce the framework of *NormalNet*. Including four parts: the generation of patch, the introduction of GNF [3], the voxelization strategy, and the proposed scheme.

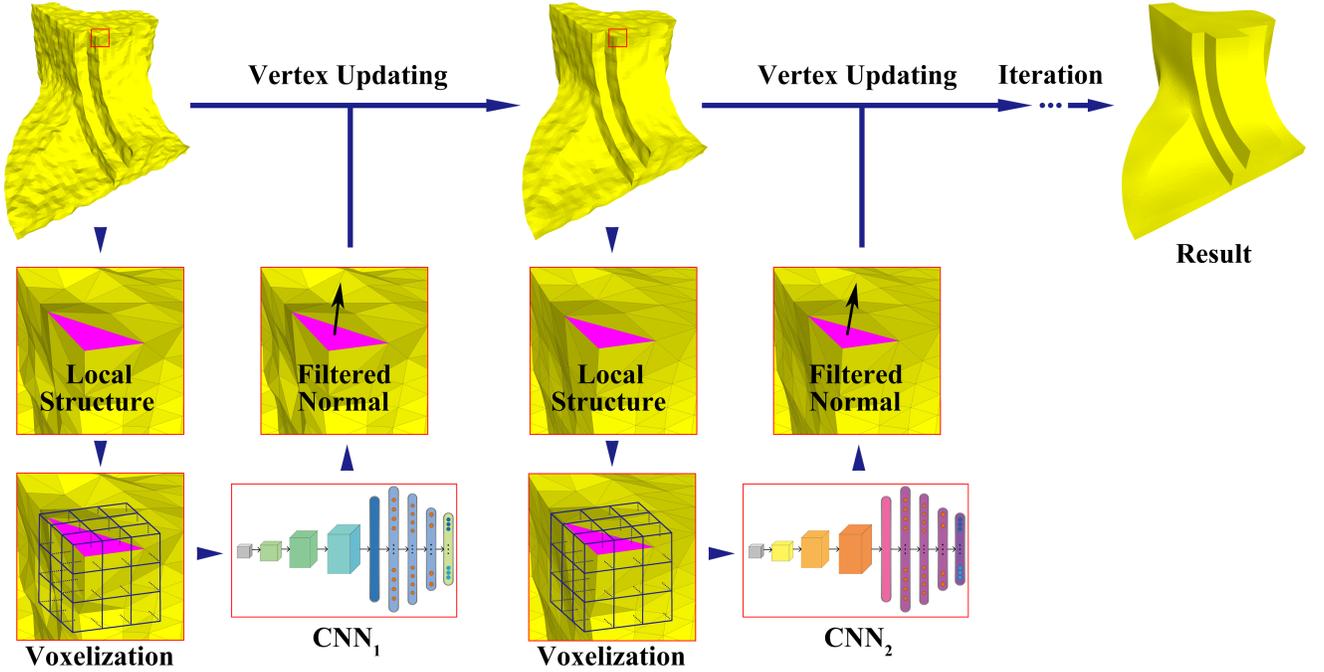


Fig. 1: The framework of *NormalNet*. Modules in all iterations share the same workflow: for a face, the irregular local 3D structure is converted via the voxelization strategy into the regular volumetric representation, which is then input into CNN to get the filtered normal. Finally, the vertex positions are updated to obtain the denoised mesh.

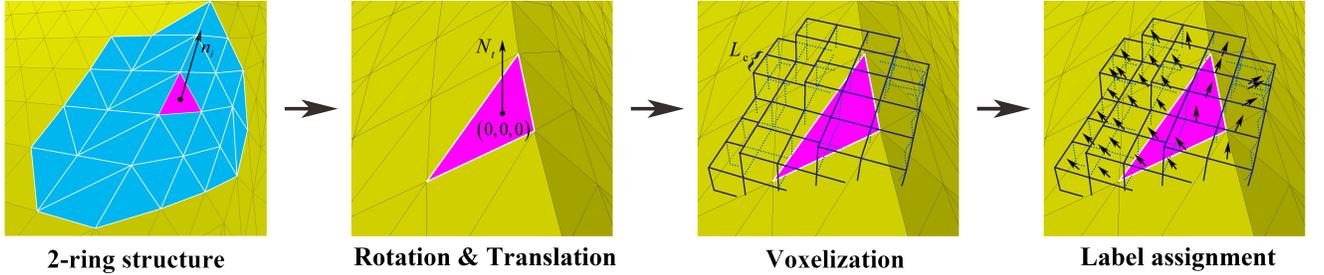


Fig. 2: Illustration of the proposed voxelization strategy. For a face in a mesh, a 2-ring patch is constructed. Two matrices that represent rotation and translation are computed for normalization. The irregular 3D structure around this face is split into small cubes. A label, which is the average normal of the faces in a cube, is then assigned to the cube.

3.1 The generation of patch

As in mesh denoising, patch is a commonly used structure, so we describe the generation of r -ring patch first. Given a face f_i as the center of patch P , an r -ring patch of f_i is generated by finding all the faces that share at least one vertex with the faces in P , and adding them into P for r times. Two examples of 1-ring and 2-ring patches are shown in Fig. 4.

3.2 The guided normal filtering

Since our scheme mimics the framework of GNF [3], we briefly introduce GNF.

GNF is an iterative scheme, in which the face normal filtering is repeated for N_f times. For a face f_i , the guided filtering is applied to obtain the denoised normal:

$$\mathbf{n}'_i = e_i \sum_{f_j \in \mathcal{N}_i} a_j G_d(c_i, c_j) G_g(\mathbf{g}_i, \mathbf{g}_j) \mathbf{n}_j \quad (1)$$

where c_j , \mathbf{n}_j and \mathbf{g}_j are the centre, face normal and guidance normal of f_j ; \mathcal{N}_i is a set of the geometrical neighbouring

faces of f_i ; e_i is a normalization factor used to ensure that \mathbf{n}'_i is a unit vector. G_d and G_g are the Gaussian kernels [45], which are computed by:

$$G_g = \exp\left(-\frac{|\mathbf{n}_i - \mathbf{n}_j|^2}{2\mu_g^2}\right), \quad (2)$$

$$G_d = \exp\left(-\frac{|c_i - c_j|^2}{2\mu_d^2}\right) \quad (3)$$

where μ_d and μ_g are the Gaussian function parameters, μ_d is usually twice the average distance between adjacent face centres, μ_g is usually different for different meshes. Following the idea of [46], after each filtering, the position updating of the vertices is repeated for N_v times to obtain the denoised mesh.

The guidance normal of f_i is generated as follows. For each face suppose that P is a 1-ring patch that contains f_i . The consistency $C(P)$ of P is calculated as [3]:

$$C(P) = D(P) \cdot R(P) \quad (4)$$

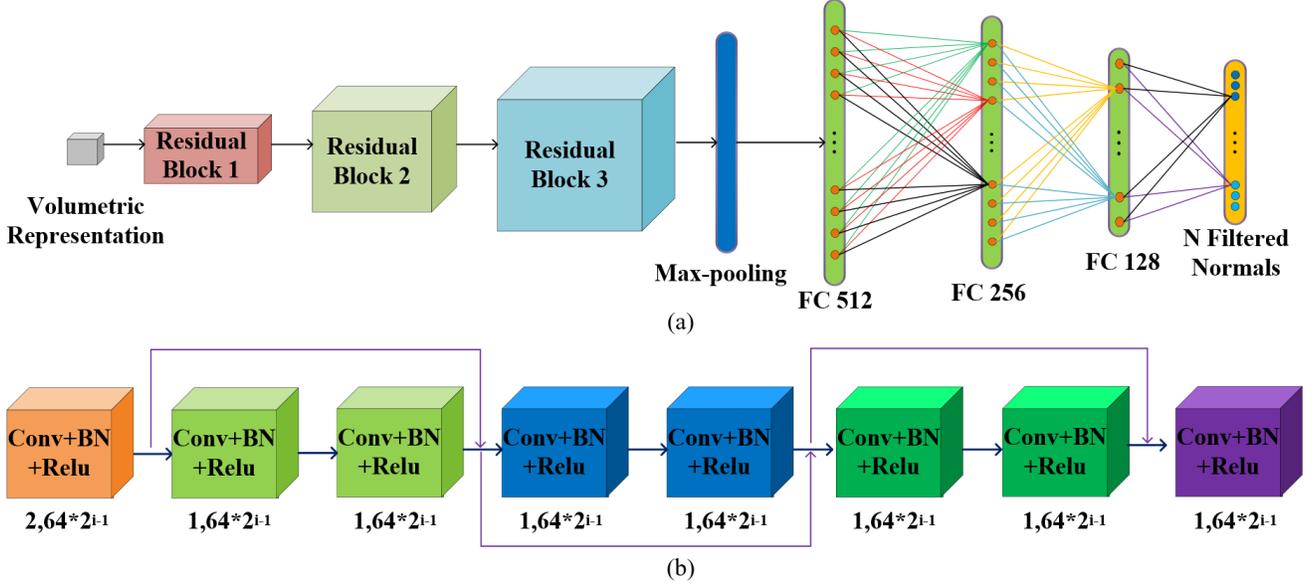


Fig. 3: (a). The architecture of the deep network in *NormalNet*. (b). The structure of residual blocks $i, [2, 64 * 2^{i-1}]$ means that the convolution stride is 2 and the channel number is $64 * 2^{i-1}$.

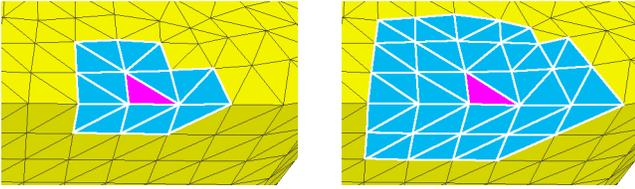


Fig. 4: Two examples of a 1-ring patch (left) and a 2-ring patch (right), f_i is colored with purple.

where $D(P)$ is the most significant face normal difference in P and calculated by:

$$D(P) = \max_{f_i, f_j \in P} |\mathbf{n}_i - \mathbf{n}_j| \quad (5)$$

where f_i and f_j represent a pair of faces within P . $R(P)$ represents the saliency of P , which is computed by the saliency of edges in P :

$$R(P) = \frac{\max_{e_i \in P} \varphi(e_i)}{\varepsilon + \sum_{e_i \in P} \varphi(e_i)} \quad (6)$$

ε is a small positive number to prevent the denominator from being zero, and $\varphi(e_i)$ is the saliency of an edge e_i :

$$\varphi(e_i) = |\mathbf{n}_{i1} - \mathbf{n}_{i2}| \quad (7)$$

where \mathbf{n}_{i1} and \mathbf{n}_{i2} are the normals of the incident faces of e_i . Finally, the most consistent patch is chosen, and the average normal of the patch is regarded as the guidance normal.

The guidance normals generated by the above method have been proven to be effective on simple structures. However, the calculation of consistency is only based on the difference between face normals, and the structure information has not been fully considered. As mentioned before, scanned meshes may contain noise with huge face normal difference. This method will not be working well for such cases.

Rather than designing a method that works well on these meshes manually, we employed CNNs to obtain the learned filtered normals.

3.3 The voxelization strategy

The key to use CNNs for mesh denoising is the transformation of the irregular local structure around a face into a regular form such that the structure information is preserved and the CNN convolution operations are easily performed.

An illustration of the proposed voxelization strategy is shown in Fig. 2. The normalization is applied to improve the robustness of the strategy first. The normalization process involves two operations: rotation and translation. In this way, all faces are normalized to a similar direction and position. Specifically, for a face f_i , a 2-ring patch is constructed. The average normal of this patch is \mathbf{n}_i . We then compute two matrices: \mathbf{W}_r , which represents the rotation from n_i to a specific angle N_t , and \mathbf{W}_t , which represents the translation from the face centre c_i to $(0,0,0)$. The whole mesh is then rotated and translated by means of \mathbf{W}_r and \mathbf{W}_t . Supposing v_i is the coordinate of a vertex i in the mesh, the new position of v_i after normalization is:

$$v'_i = \mathbf{W}_t \mathbf{W}_r v_i \quad (8)$$

After normalization, the space of the local mesh structure around f_i is split into regular cubes denoted by $\{B_{x,y,z} | x, y, z \in [-T_s, T_s]\}$, where T_s is the parameter that determines the number of cubes and $B_{0,0,0}$ is located at the origin. The rest issue is to determine the size of each cube. In our work, the side length L_c of the cubes is computed as:

$$L_c = \frac{d_s}{\alpha_c} \quad (9)$$

where d_s is the average distance between adjacent faces in the noisy mesh and α_c is the parameter that controls the size of the cubes.

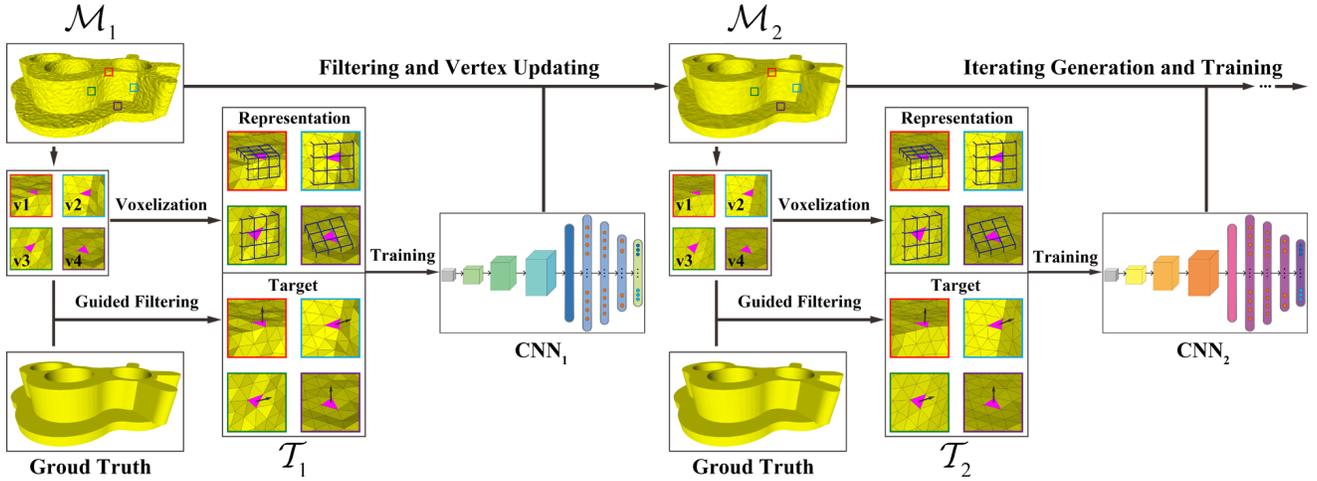


Fig. 5: The framework of generating training sets and training.

For each cube, we employ the fast 3D triangle-box overlap testing strategy [47] to find faces that overlap with this cube. If at least one face overlaps with this cube, the label of this cube is assigned as the average normal of all the overlapped faces, denoted by \mathcal{B} ; otherwise, the label is set to $(0,0,0)$. In this way, we convert the irregular local mesh structure into the regular volumetric representation $\mathcal{V} = [\mathcal{B} \in \mathbb{R}^{(2T_s+1)^3}]$. \mathcal{V} is then used as the input of the network.

In our experiment, we set $T_s = 20$, $\alpha_c = 8$ and $N_t = (0, 1, 0)$. Under these conditions, \mathcal{V} is a $41 \times 41 \times 41 \times 3$ matrix that contains the most 3-ring structure around f_i . Each face is split into about 40~60 cubes, which is sufficient to represent the shape information. Smaller T_s and α_c reduce the amount of information in \mathcal{V} and lead to unsatisfactory results, whereas larger parameters can improve the performance slightly but greatly increase the training time.

3.4 The proposed scheme

The proposed scheme is also an iterative scheme which is repeated for N_f times. During each iteration, for a face f_i in a mesh, the voxelization strategy is employed to transform the irregular local mesh structure around f_i into the regular volumetric representation. Then a CNN takes the volumetric representation as input, and outputs the filtered normals. Since the value of μ_g in GNF greatly affects the denoising results and is often different for different meshes. Therefore the output of the network contains N filtered normals with different μ_g . At last the positions of the vertices are updated according to the selected filtered normals by N_v times.

The network architecture is shown in Fig. 3. It contains three residual blocks, a global max-pooling layer and four fully connected layers. The numbers of channels of the residual blocks are 64, 128 and 256. All the convolution layers use $3 * 3 * 3$ filters except the first layer, which uses $5 * 5 * 5$ filters. Down-sampling is performed by a convolution operation with a stride of 2 in the first layer of each residual block. The network ends with four fully connected layers: the first three have 512, 256 and 128 channels. The fourth aims to predict the three coordinates of N filtered normals and thus contains $3 * N$ channels. All layers are equipped

CNN _i	1	2	3	4	5	6
Iteration numbers	[1,1]	[2,2]	[3,3]	[4,5]	[6,10]	[11, N _f]

TABLE 1: The settings of the corresponding iteration numbers for each CNN_i.

with batch normalization and ReLU, except the last layer is equipped with Tanh to ensure the output lies in $[-1,1]$.

The network architecture is inspired by the philosophy of ResNet [48] and VGGNet [49]. The purpose of *NormalNet* is to estimate accurate filtered normals from the noisy signal. However, as the network goes deeper, abundant information beneficial to filtering normals from the input can vanish or "wash out" by the time it reaches the output layer. To address this problem, we adopt the shortcut connection from ResNet to directly pass the early feature map to the later layers. This greatly increases the forward flow of information and thus contributes to the prediction of face normals. In addition, during the backpropagation process, a shortcut path adds an extra component to the gradients compared to the plain network, which can mitigate the vanishing gradient problem, thereby accelerating the training process.

In our experiment, we set $N = 6$, and the output of the CNN contains the filtering results of $\mu_g = 0.25, 0.3, 0.35, 0.4, 0.45, 0.5$.

4 NORMALNET TRAINING

In this section, we introduce the training of *NormalNet*. Including two parts: the iterative training and the training details.

4.1 The iterative training

The process of generating training sets and training is illustrated in Fig. 5. For each CNN_i, a specified training data set \mathcal{T}_i is generated from a group of meshes named by \mathcal{M}_i and the corresponding ground truth. \mathcal{T}_i is composed of numerous training tuples, each of which consists of a volumetric representation and N target normals. For a face f_i from a mesh in \mathcal{M}_i , the volumetric representation is

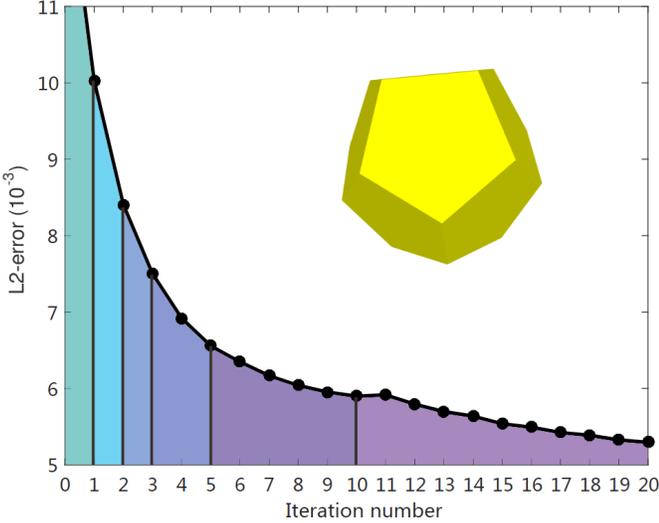


Fig. 6: Illustration of the L2 error results on the model *Twelve*; each colour represents a CNN.

Para	<i>Fandisk</i>	<i>Table</i>	<i>Joint</i>	<i>Twelve</i>	<i>Block</i>
N_f	10	15	5	25	20
N_v	20	15	15	10	30
μ_g	0.25	0.4	0.25	0.3	0.3
Para	<i>Bunny</i>	<i>Angel</i>	<i>Iron</i>	<i>Pierrot</i>	<i>Rocker-arm</i>
N_f	2	3	10	10	10
N_v	5	4	10	10	10
μ_g	0.3	0.3	0.35	0.35	0.25
Para	<i>Eagle</i>	<i>Gargoyle</i>	<i>Balljoint</i>	<i>Boy01F</i>	<i>Boy02F</i>
N_f	4	5	4	14	7
N_v	5	10	10	20	20
μ_g	0.4	0.3	0.4	0.4	0.35
Para	<i>Cone04V1</i>	<i>Girl02V1</i>	<i>Cone16V2</i>	<i>Girl01V2</i>	-
N_f	20	15	10	3	-
N_v	20	20	10	15	-
μ_g	0.45	0.45	0.3	0.4	-

TABLE 2: The settings of N_f , N_v and μ_g .

obtained by applying the voxelization strategy on f_i . The target normals are obtained by employing GNF, whereas the ground truth normals are adopted as the guidance normals.

$$\mathbf{n}'_i = e_i \sum_{f_j \in \mathcal{N}_c} a_j G_d(c_i, c_j) G_g(\mathbf{gn}_i, \mathbf{gn}_j) \mathbf{n}_j \quad (10)$$

where \mathbf{gn}_i and \mathbf{gn}_j are the ground truth normals of f_i and f_j , The other parameters are the same as defined in Eq.(1).

To make the training process balance with respect to various features. Suppose the maximum angle difference in the 2-ring patch of a face is A_p . All the faces in \mathcal{M}_i are divided into 4 categories and we randomly select the same number of faces in each category for training:

- v1: $A_p > 80^\circ$, large edge region.
- v2: $50^\circ < A_p \leq 80^\circ$, small edge region.
- v3: $20^\circ < A_p \leq 50^\circ$, curved region.
- v4: $A_p \leq 20^\circ$, smooth region.

Initially, \mathcal{M}_1 is composed of noisy meshes that their ground truth are already known and without any processing. When $i > 1$, \mathcal{M}_i will be obtained by applying filtering on \mathcal{M}_{i-1} , which is performed by CNN_{i-1} , the parameters

used in filtering are $\mu_g = 0.4$ and $N_v = 20$. The generation of the training data and the network training are alternately performed iteratively.

4.2 Training details

The loss function is defined as the MSE between N output normals and the target normals. We use the truncated normal distribution to initialize the weights and train the network from scratch. For the optimization method, we choose the Adam algorithm with a mini-batch size of 80, and the parameters for the Adam optimizer are $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 1e - 8$, which are the default settings in TensorFlow. The learning rate starts at 0.0001 and decays exponentially after 5000 training steps, for which the decay rate is 0.96. Each CNN_i is trained individually. A test set that randomly selects some faces from test models is built for evaluation. The evaluation metric for the network is defined as the average angular error over the entire test set. Each network is trained for 10 epochs, and the average angular error is 1-3 degrees after 10 epochs. The network with the smallest error is selected for utilization.

In our experiment, we select 45000 faces in each category; thus, the total size of \mathcal{T}_i is 180000. The training process is executed on a computer with an Intel Core i7-7700 CPU and NVIDIA GTX1080, and each epoch is approximately 3 hours. Increasing the number of channels, the number of layers or the size of \mathcal{T}_i will not substantially improve the performance of the networks and only multiplies the training time. Halving the numbers of channels or the size of \mathcal{T}_i will also halve the training time; however, the average angular error will increase to 4-6 degrees.

5 EXPERIMENTAL RESULTS

In this section, the extensive experimental results are presented to demonstrate the performance of *NormalNet*.

5.1 Comparison study

We perform the experimental comparisons on 19 test models, including 6 synthetic models: *Joint*, *Twelve*, *Bunny*, *Fandisk*, *Table*, *Block*; 4 scanned models collected from Internet where the type of scanner is unknown: *Angel*, *Iron*, *Rocketarm*, *Pierrot*; 6 scanned models which have rich features and generated by Microsoft Kinect v1, Microsoft Kinect v2 and Microsoft Kinect v1 via the Kinect-Fusion technique [43], respectively: *Core04V1*, *Girl02V1*, *Core16V2*, *Girl01V2*, *Boy01F*, *Boy02F*; and 3 scanned models generated by laser scanners [50]: *Eagle*, *Gargoyle* and *Balljoint*. For the synthetic models, the noise type in *Fandisk*, *Table*, *Bunny* and *Block* is Gaussian white noise, while that of *Joint* and *Twelve* is impulsive noise.

We compare *NormalNet* with several state-of-the-art algorithms in terms of objective and subjective evaluations. The compared algorithms are 1) guided normal filtering (GNF) [3], 2) L0 minimization optimization (LOM) [10], 3) BI-normal filtering (BI) [2], 4) cascaded normal regression (CNR) [43], 5) graph-based normal filtering (GGNF) [7], and 6) normal-voting-tensor-based scheme (VT) [50]. The source codes of GNF, LOM, BI, CNR and GGNF are kindly provided by their authors or implemented by a third party, while the author of VT provides the input models and their denoising results.

model	Noise Level	<i>NormalNet</i>	Metrics	L0M [10]	BI [2]	GNF [3]	CNR [43]	GGNF [7]	<i>NormalNet</i>
<i>Fandisk</i>	0.3	CS_1	$E_v (\times 10^{-3})$	1.850	1.509	1.458	1.564	1.430	1.281
			E_a	10.141	11.670	7.615	4.653	6.130	4.560
<i>Table</i>	0.3	CS_1	$E_v (\times 10^{-3})$	1.961	1.571	1.894	1.669	1.378	1.372
			E_a	12.348	18.635	17.544	18.912	15.184	18.810
<i>Joint</i>	0.2	CS_1	$E_v (\times 10^{-4})$	2.429	1.780	1.428	1.434	1.438	1.403
			E_a	11.181	5.489	5.920	4.390	2.857	2.625
<i>Twelve</i>	0.5	CS_1	$E_v (\times 10^{-3})$	20.00	11.68	5.955	7.427	5.285	5.132
			E_a	12.147	20.038	11.099	5.734	8.550	5.290
<i>Block</i>	0.4	CS_1	$E_v (\times 10^{-3})$	9.273	4.895	5.417	5.944	5.131	5.331
			E_a	10.722	15.689	10.438	6.725	10.007	5.748
<i>Bunny</i>	0.2	CS_1	$E_v (\times 10^{-6})$	7.897	7.727	7.713	7.879	7.673	7.660
			E_a	9.359	11.008	7.494	7.649	7.246	6.963
<i>Boy01F</i>	Scanned	CS_1	$E_v (\times 10^{-4})$	8.119	8.120	8.170	8.179	8.106	8.098
			E_a	19.182	20.181	16.592	16.903	16.266	15.994
<i>Boy02F</i>	Scanned	CS_1	$E_v (\times 10^{-3})$	8.446	8.514	8.398	8.392	8.302	8.347
			E_a	16.601	17.056	14.491	14.931	14.445	13.966
<i>Cone04V1</i>	Scanned	CV_1	$E_v (\times 10^{-3})$	3.569	2.781	2.657	2.806	2.575	2.568
			E_a	37.618	22.144	15.658	15.670	15.836	13.157
<i>Girl02V1</i>	Scanned	CV_1	$E_v (\times 10^{-3})$	1.934	1.899	1.751	1.658	1.769	1.634
			E_a	37.826	26.353	19.707	20.121	19.672	17.903
<i>Cone16V2</i>	Scanned	CV_2	$E_v (\times 10^{-3})$	14.539	16.508	8.998	8.948	8.690	8.642
			E_a	27.872	12.642	10.805	8.468	9.862	8.731
<i>Girl01V2</i>	Scanned	CV_2	$E_v (\times 10^{-3})$	5.461	5.261	5.261	5.425	5.226	5.171
			E_a	28.401	18.098	18.098	14.627	18.487	14.017
Average	-	-	E_v	7.123	6.020	4.925	5.110	4.750	4.719
			E_a	19.449	16.583	12.955	11.565	12.045	10.647

TABLE 3: Performance comparisons between *NormalNet* and the state-of-the-art methods.

5.2 Parameter settings

As shown in Fig. 6, during the denoising process for most meshes, the L2-error decreases rapidly during the first three iterations and decreases slowly after ten iterations. In order to design a lightweight network, the iteration numbers are divided into six intervals, each of which corresponds to a specific CNN_i , as listed in Table 1. Thus, the training cost decreases by more than 70% at the price of slightly decreased performance of CNN, the average angular error will increase 0.1-0.15 degree.

Three *NormalNet*, namely, CV_1 , CV_2 and CS_1 , are trained on Kinect-v1 training set (73 meshes), Kinect-v2 training set (73 meshes) and a remake of synthetic training set (60 meshes, where some meshes are excluded from the training sets for experiments) provided by [43]. The test models *Cone04V1*, *Girl02V1*, *Cone16V2*, and *Girl01V2* are denoised by the corresponding networks CV_1 and CV_2 , and all the other test models are denoised by CS_1 . The settings of the parameters N_f , N_v and μ_g and the parameters in other schemes refer to the settings used in [3] and [7]. The parameter settings of N_f , N_v and μ_g are shown in Table 2.

5.3 Objective performance comparison

Two error metrics [20] are employed to evaluate the objective denoising results of the models which have the ground truth:

- E_a : the mean angle square error, which represents the accuracy of the face normal;
- E_v : the L2 vertex-based mesh-to-mesh error, which represents the accuracy of a vertex’s position.

We compare the objective performance on 12 models. The comparison results of E_a and E_v are shown in Table 3, where the best results are bolded, *NormalNet* performs best for 10 models on E_a and 10 models on E_v , which achieves the best performance with respect to both metrics on most test models. CNR achieves the second best average results on E_a , which proves CNR is superior in estimating face normals. However, GNF and GGNF achieve better average results than CNR on E_v , which proves that filtering-based schemes perform better in recovering vertex positions. *NormalNet* achieves the best average results on E_a and E_v .

5.4 Subjective Performance Comparison

5.4.1 Results on synthetic models

The subjective performance comparison results of six synthetic models are illustrated in Figs. 7, 8 and 9.

Fig. 7 presents the denoising results of two models with curved surfaces. The zoomed-in view illustrates that our scheme introduces fewer pseudo-features than other schemes. In Fig. 8, our scheme achieves similar performance to that of GGNF in these feature regions. The corner is recovered well, and the edge is sharp and clean. In *Block*, the highlighted region in the red window has a higher triangulation density. Benefiting from the voxelization strategy, our scheme can preserve the structure information well and is thus less sensitive to the sampling irregularity. In Fig. 9, we perform a comparison on synthetic meshes with impulsive noise. In *Table*, both our scheme and GGNF produce the best feature recovery results. In *Joint*, the edge length of our scheme is closest to the ground truth.

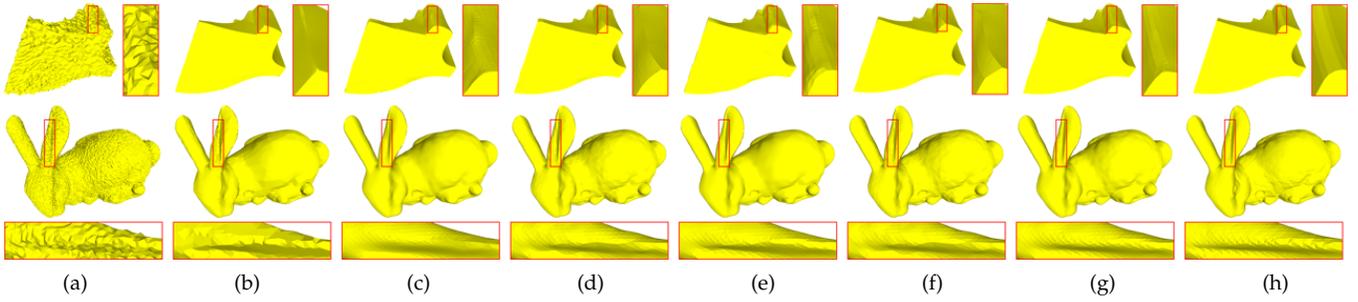


Fig. 7: Illustration of the denoising results on the models *Fandisk* and *Bunny*; the zoomed-in view of *Bunny* has been rotated. (a) to (h) are the noisy mesh; the results of LOM [10], BI [2], GNF [3], CNR [43], GGNF [7] and *NormalNet*; and the ground truth.

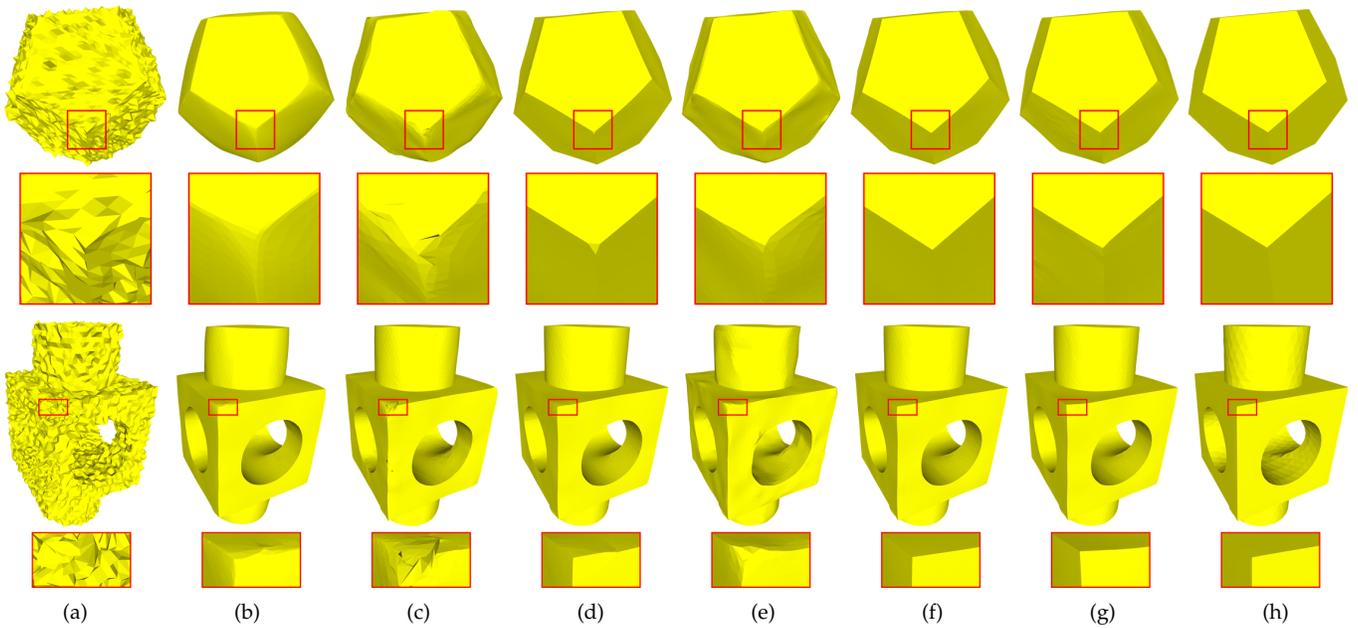


Fig. 8: Illustration of the denoising results on the models *Twelve* and *Block*. (a) to (h) are the noisy mesh; the results of LOM [10], BI [2], GNF [3], CNR [43], GGNF [7] and *NormalNet*; and the ground truth.

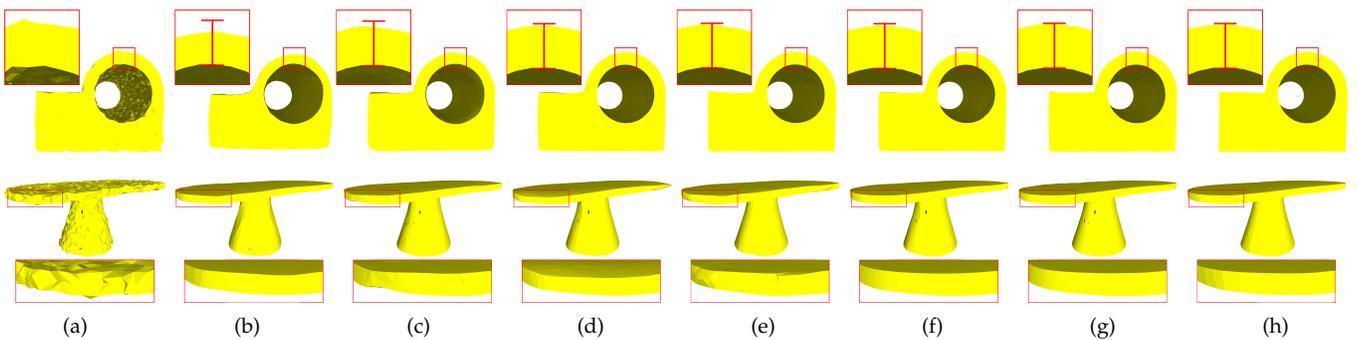


Fig. 9: Illustration of the denoising results on the models *Joint* and *Table*. The red line in *Joint* is the length of the ground truth. (a) to (h) are the noisy mesh; the results of LOM [10], BI [2], GNF [3], CNR [43], GGNF [7] and *NormalNet*; and the ground truth.

5.4.2 Results on scanned models

We further provide the comparison results for models with different scanners. As illustrated in Figs. 10 and 11, our scheme preforms well on models where the type of scanner is unknown. For fair comparison, CNR is also trained on the synthetic training set. In Fig. 10, for *Iron*, both our scheme and CNR introduce fewer pseudo-features than the other schemes. However, in *Rocketarm* and *Angel*, CNR over-smooths the edges in the red boxes, whereas our scheme still produces satisfactory results. In Fig. 11, for *Pierrot*, the region in the red box is corrupted by serrated noise. For GNF and GGNF, accurate guidance normals are difficult to compute under this type of noise. Thus, the denoising result is corrupted by pseudo-features. Furthermore, CNR succeeds in removing the serrated noise but fails to recover the edges around the eyes in the red box. Our scheme finds a balance between introducing pseudo-features and over-smoothing. The codes of LOM and BI could not process this region.

In Fig. 12, we compare *NormalNet* with (VT) [50] on the models provided by the authors, which are generated by laser scanners. Our scheme produces better feature recovery results than VT on all three models that contain complex structures, which further verifies the capability of *NormalNet*.

In Figs. 13 and 14, the models are generated by Microsoft Kinect V1 and V2 and provided by the author of CNR. However, we do not have sufficient data to train *NormalNet* for the models generated by Microsoft Kinect v1 via the Kinect-Fusion technique. Therefore, CS_1 is employed to denoise these models. In Fig. 13, our scheme outputs similar denoising results as CNR and GGNF. In Fig. 14, our scheme achieves the best smoothing result and the other schemes fail to remove noise in *Cone04V1*. In *Girl02V1* and *Girl01V2*, both CNR and our scheme avoid introducing pseudo-features. In *Cone16V2*, most schemes achieve similar feature recovery results.

6 CONCLUSION

In this paper, we present a learning-based normal filtering scheme for mesh denoising. The scheme maps the guided normal filtering into a deep network and follows the iterative framework of filtering-based scheme. During each iteration, first, to facilitate the 3D convolution operations, the voxelization strategy is applied on each face in a mesh to transform the irregular local structure into the regular volumetric representation. Second, instead of the guidance normal generation and the guided filtering in GNF, the output of voxelization is then input into a CNN to estimate accurate filtered normals. Finally, the vertex positions are updated according to the filtered normals. What's more, the iterative training framework is proposed for effectively training. The experimental results show that our scheme outperforms state-of-the-art works with respect to both objective and subjective quality metrics and can effectively remove noise while preserving the original features and avoiding pseudo-features.

REFERENCES

- [1] H. Yagou, Y. Ohtake, and A. Belyaev, "Mesh smoothing via mean and median filtering applied to face normals," in *Geometric Modeling and Processing*, 2002, pp. 124–131.
- [2] M. Wei, J. Yu, W.-M. Pang, J. Wang, J. Qin, L. Liu, and P.-A. Heng, "Bi-normal filtering for mesh denoising," *IEEE Transactions on Visualization and Computer Graphics*, vol. 21, no. 1, pp. 43–55, 2015.
- [3] W. Zhang, B. Deng, J. Zhang, S. Bouaziz, and L. Liu, "Guided mesh normal filtering," in *Computer Graphics Forum*, vol. 34, no. 7. Wiley Online Library, 2015, pp. 23–34.
- [4] W. Zhao, X. Liu, S. Wang, and D. Zhao, "Multi-scale similarity enhanced guided normal filtering," in *Advances in Multimedia Information Processing – PCM 2017*. Springer International Publishing, 2018, pp. 645–653.
- [5] R. Wang, W. Zhao, S. Liu, D. Zhao, and C. Liu, "Feature-preserving mesh denoising based on guided normal filtering," in *Advances in Multimedia Information Processing – PCM 2017*. Springer International Publishing, 2018, pp. 920–927.
- [6] T. Li, J. Wang, H. Liu, and L.-g. Liu, "Efficient mesh denoising via robust normal filtering and alternate vertex updating," *Frontiers of Information Technology and Electronic Engineering*, vol. 18, no. 11, pp. 1828–1842, 2017.
- [7] W. Zhao, X. Liu, S. Wang, X. Fan, and D. Zhao, "Graph-based feature-preserving mesh normal filtering," *IEEE Transactions on Visualization and Computer Graphics*, 2019.
- [8] X. Lu, Z. Deng, and W. Chen, "A robust scheme for feature-preserving mesh denoising," *IEEE Transactions on Visualization and Computer Graphics*, vol. 22, no. 3, pp. 1181–1194, 2016.
- [9] M. Wei, L. Liang, W.-M. Pang, J. Wang, W. Li, and H. Wu, "Tensor voting guided mesh denoising," *IEEE Transactions on Automation Science and Engineering*, vol. 14, no. 2, pp. 931–945, 2017.
- [10] L. He and S. Schaefer, "Mesh denoising via l0 minimization," *ACM Transactions on Graphics*, vol. 32, no. 4, p. 64, 2013.
- [11] R. Wang, Z. Yang, L. Liu, J. Deng, and F. Chen, "Decoupling noise and features via weighted l1-analysis compressed sensing," *ACM Transactions on Graphics*, vol. 33, no. 2, p. 18, 2014.
- [12] S. Yoshizawa, A. Belyaev, and H.-P. Seidel, "Smoothing by example: Mesh denoising by averaging with similarity-based weights," in *IEEE International Conference on Shape Modeling and Applications*. IEEE, 2006, pp. 9–9.
- [13] G. Rosman, A. Dubrovina, and R. Kimmel, "Patch-collaborative spectral point-cloud denoising," in *Computer Graphics Forum*, vol. 32, no. 8. Wiley Online Library, 2013, pp. 1–12.
- [14] J. Digne, "Similarity based filtering of point clouds," in *Computer Vision and Pattern Recognition Workshops*. IEEE, 2012, pp. 73–79.
- [15] Y. Chen and T. Pock, "Trainable nonlinear reaction diffusion: A flexible framework for fast and effective image restoration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1256–1272, 2015.
- [16] K. Zhang, W. Zuo, and L. Zhang, "Ffdnet: Toward a fast and flexible solution for cnn based image denoising," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4608–4622, 2017.
- [17] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [18] S. Fleishman, I. Drori, and D. Cohen-Or, "Bilateral mesh denoising," in *ACM Transactions on Graphics*, vol. 22, no. 3, 2003, pp. 950–953.
- [19] T. R. Jones, F. Durand, and M. Desbrun, "Non-iterative, feature-preserving mesh smoothing," in *ACM Transactions on Graphics*, vol. 22, no. 3, 2003, pp. 943–949.
- [20] Y. Zheng, H. Fu, O. K.-C. Au, and C.-L. Tai, "Bilateral normal filtering for mesh denoising," *IEEE Transactions on Visualization and Computer Graphics*, vol. 17, no. 10, pp. 1521–1530, 2011.
- [21] G. Petschnigg, R. Szeliski, M. Agrawala, M. Cohen, H. Hoppe, and K. Toyama, "Digital photography with flash and no-flash image pairs," in *ACM SIGGRAPH 2004 Papers*, 2004, pp. 664–672.
- [22] H. Zhang, C. Wu, J. Zhang, and J. Deng, "Variational mesh denoising using total variation and piecewise constant function space," *IEEE Transactions on Visualization and Computer Graphics*, vol. 21, no. 7, pp. 873–886, 2015.
- [23] "Robust mesh denoising via vertex pre-filtering and l1-median normal filtering," *Computer Aided Geometric Design*, vol. 54, pp. 49 – 60, 2017.

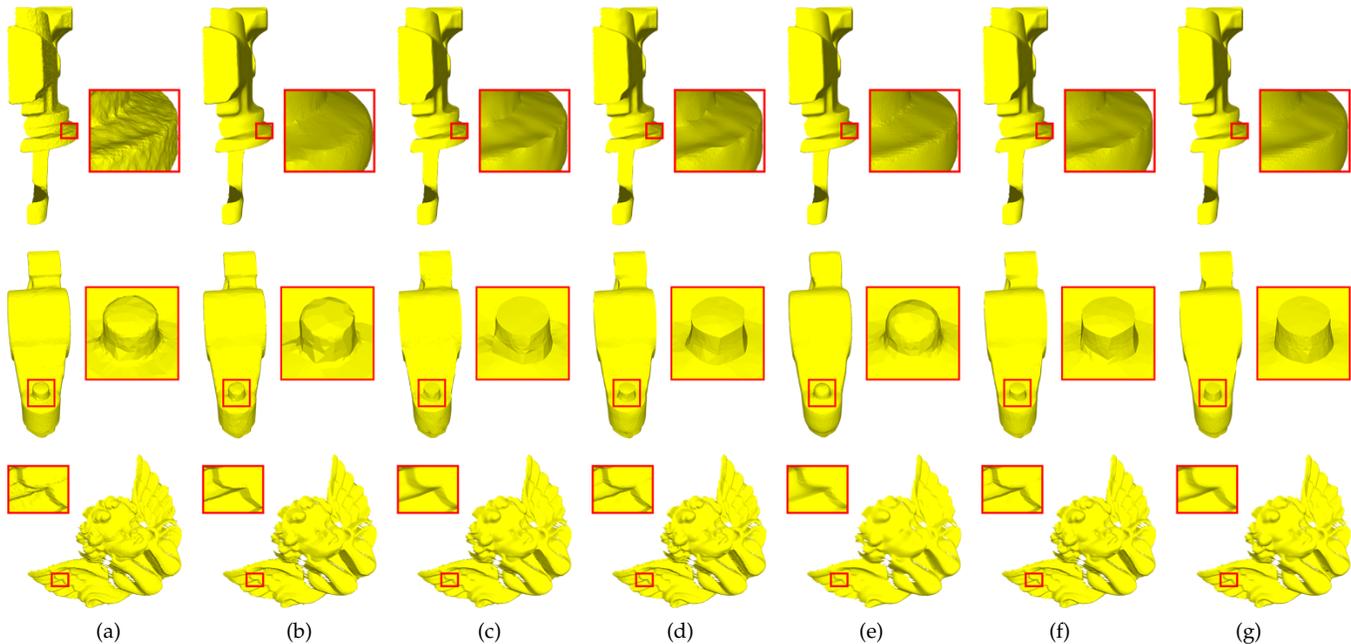


Fig. 10: Illustration of the denoising results on the models *Iron*, *Rocketarm* and *Angel*, which are generated by unknown scanners. (a) to (g) are the noisy mesh and the results of LOM [10], BI [2], GNF [3], CNR [43], GGNF [7] and *NormalNet*.

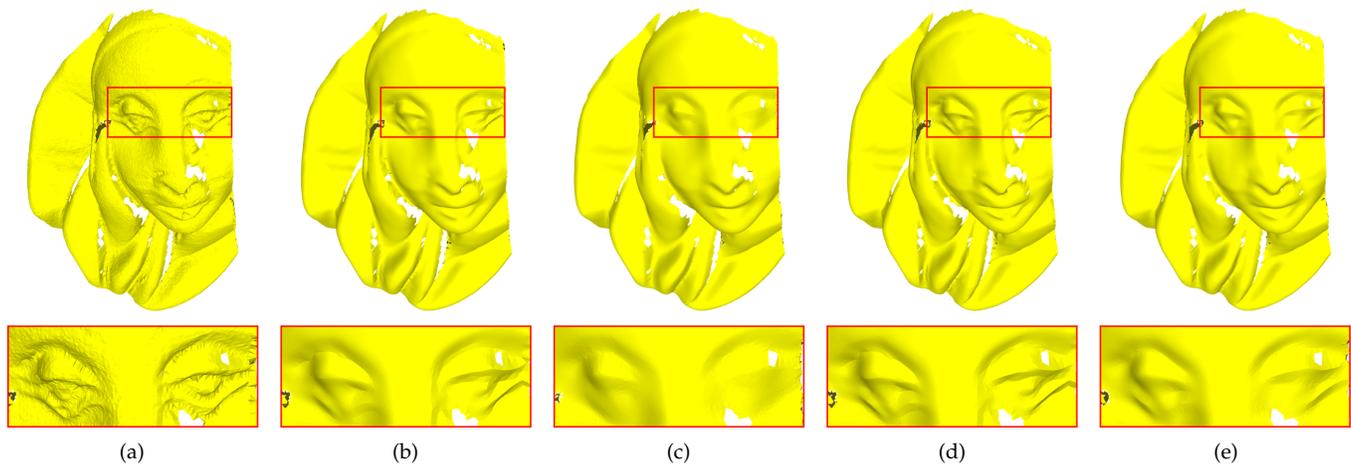


Fig. 11: Illustration of the denoising results on the scanned model *Pierrot*, which is generated by unknown scanners. (a) to (e) are the noisy mesh and the results of GNF [3], CNR [43], GGNF [7] and *NormalNet*.

- [24] S. K. Yadav, U. Reitebuch, and K. Polthier, "Robust and high fidelity mesh denoising," *IEEE Transactions on Visualization and Computer Graphics*, vol. 25, no. 6, pp. 2304–2310, 2019.
- [25] H. Su, S. Maji, E. Kalogerakis, and E. Learned-Miller, "Multi-view convolutional neural networks for 3d shape recognition," in *IEEE International Conference on Computer Vision*, Dec 2015, pp. 945–953.
- [26] B. Shi, S. Bai, Z. Zhou, and X. Bai, "Deeppano: Deep panoramic representation for 3-d shape recognition," *IEEE Signal Processing Letters*, vol. 22, no. 12, pp. 2339–2343, Dec 2015.
- [27] Z. Wu, S. Song, A. Khosla, F. Yu, L. Zhang, X. Tang, and J. Xiao, "3d shapenets: A deep representation for volumetric shapes," in *The IEEE Conference on Computer Vision and Pattern Recognition*, June 2015.
- [28] D. Maturana and S. Scherer, "Voxnet: A 3d convolutional neural network for real-time object recognition," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*, Sept 2015, pp. 922–928.
- [29] P. Wang, Y. Liu, Y. Guo, C. Sun, and X. Tong, "O-CNN: octree-based convolutional neural networks for 3d shape analysis," *CoRR*, vol. abs/1712.01537, 2017.
- [30] X. Han, Z. Li, H. Huang, E. Kalogerakis, and Y. Yu, "High-resolution shape completion using deep neural networks for global structure and local geometry inference," *CoRR*, vol. abs/1709.07599, 2017.
- [31] Q. Tan, L. Gao, Y. Lai, J. Yang, and S. Xia, "Mesh-based autoencoders for localized deformation component analysis," *CoRR*, vol. abs/1709.04304, 2017.
- [32] B. Davide, M. Jonathan, R. Emanuele, B. M. M., and C. Daniel, "Anisotropic diffusion descriptors," *Computer Graphics Forum*, 2016.
- [33] L. Yi, H. Su, X. Guo, and L. J. Guibas, "Syncspecnn: Synchronized spectral cnn for 3d shape segmentation," *CoRR*, vol. abs/1612.00606, 2016.
- [34] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, "Dynamic graph cnn for learning on point clouds," *CoRR*, vol. abs/1801.07829, 2018.
- [35] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," *CoRR*, vol. abs/1612.00593, 2016.
- [36] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," in

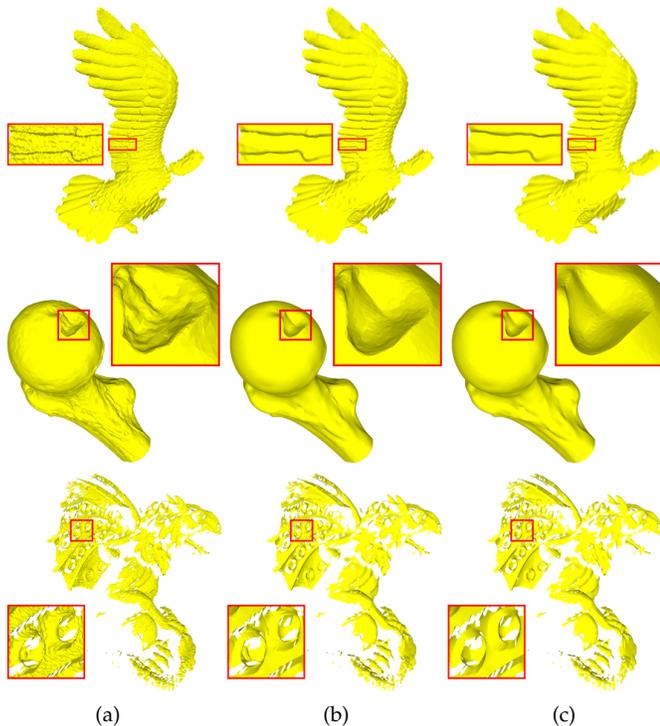


Fig. 12: Illustration of the denoising results of the scanned models *Eagle*, *Gargoyle* and *Balljoint*, which are generated by laser scanners. (a) to (c) are the noisy mesh and the results of VT [3] and *NormalNet*.

- Advances in Neural Information Processing Systems*, 2017, pp. 5099–5108.
- [37] Y. Wang, Y. Sun, Z. Liu, S. E. Sarma, M. M. Bronstein, and J. M. Solomon, “Dynamic graph cnn for learning on point clouds,” *arXiv preprint arXiv:1801.07829*, 2018.
- [38] R. Klokov and V. S. Lempitsky, “Escape from cells: Deep kd-networks for the recognition of 3d point cloud models,” in *IEEE International Conference on Computer Vision*, 2017, pp. 863–872.
- [39] G. Riegler, A. O. Ulusoy, and A. Geiger, “Octnet: Learning deep 3d representations at high resolutions,” *CoRR*, vol. abs/1611.05009, 2016.
- [40] A. Boulch and R. Marlet, “Deep learning for robust normal estimation in unstructured point clouds,” in *Computer Graphics Forum*, vol. 35, no. 5. Wiley Online Library, 2016, pp. 281–290.
- [41] R. Roveri, A. C. Öztireli, I. Pandele, and M. Gross, “Pointpronets: Consolidation of point clouds with convolutional neural networks,” *Computer Graphics Forum (Proc. Eurographics)*, vol. 37, no. 2, 2018.
- [42] M.-J. Rakotosaona, V. La Barbera, P. Guerrero, N. J. Mitra, and M. Ovsjanikov, “Pointcleannet: Learning to denoise and remove outliers from dense point clouds,” *Computer Graphics Forum*, 2019.
- [43] P.-S. Wang, Y. Liu, and X. Tong, “Mesh denoising via cascaded normal regression,” *ACM Transactions on Graphics (SIGGRAPH Asia)*, vol. 35, no. 6, 2016.
- [44] R. Hanocka, A. Hertz, N. Fish, R. Giryas, S. Fleishman, and D. Cohen-Or, “Meshcnn: A network with an edge,” *CoRR*, vol. abs/1809.05910, 2018.
- [45] C. Tomasi and R. Manduchi, “Bilateral filtering for gray and color images,” in *Sixth International Conference on Computer Vision*. IEEE, 1998, pp. 839–846.
- [46] X. Sun, P. Rosin, R. Martin, and F. Langbein, “Fast and effective feature-preserving mesh denoising,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 13, no. 5, 2007.
- [47] T. Akenine-Möller, “Fast 3d triangle-box overlap testing,” *Journal of Graphics Tools*, vol. 6, no. 1, pp. 29–33, 2001.
- [48] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.

- [49] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014.
- [50] S. K. Yadav, U. Reitebuch, and K. Polthier, “Mesh denoising based on normal voting tensor and binary optimization,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 24, no. 8, pp. 2366–2379, Aug 2018.

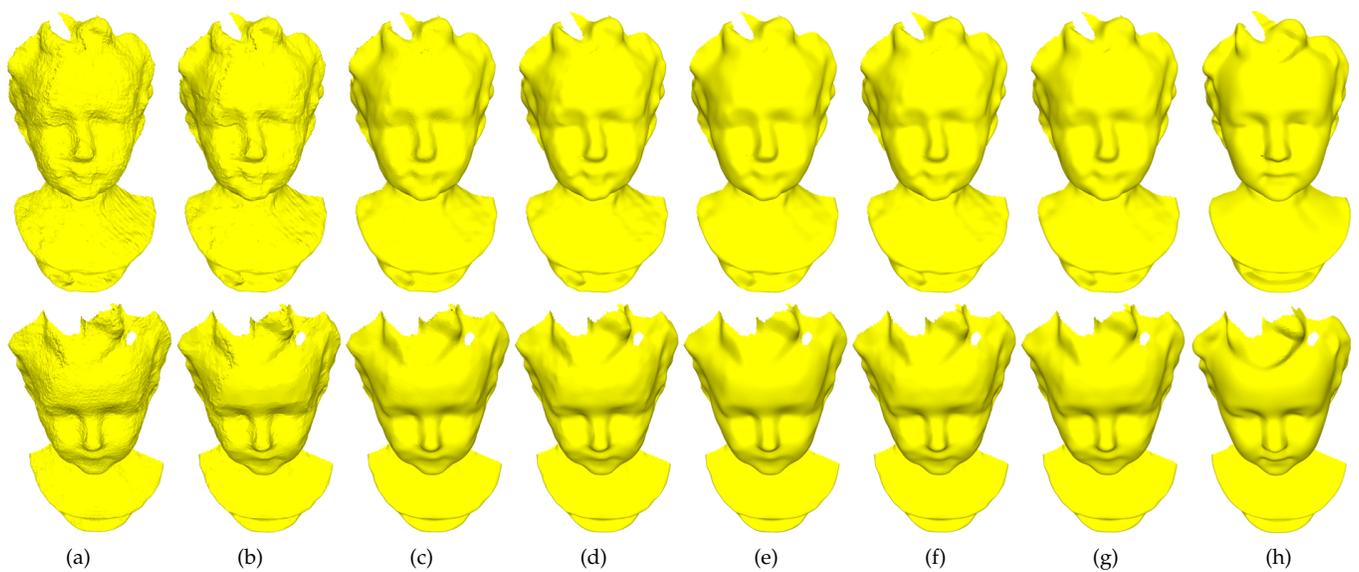


Fig. 13: Illustration of the denoising results on the models *Boy01F* and *Boy02F*, which are generated by Microsoft Kinect v1 via the Kinect-Fusion technique. (a) to (h) are the noisy mesh; the results of LOM [10], BI [2], GNF [3], CNR [43], GGNF [7] and *NormalNet*; and the ground truth.

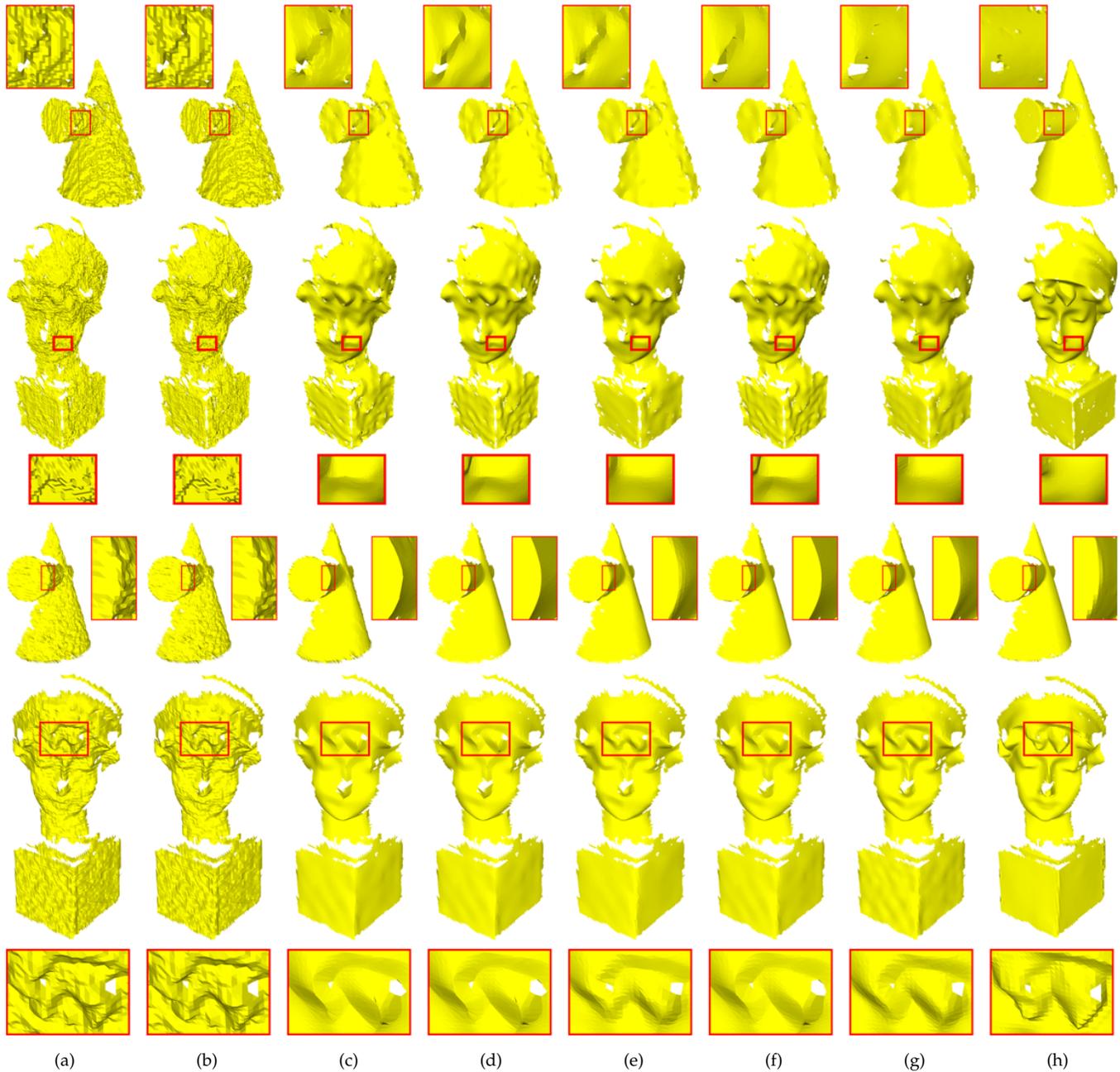


Fig. 14: Illustration of the denoising results on the models *Cone04V1*, *Girl02V1*, *Cone16V2* and *Girl01V2*, which are generated by Microsoft Kinect v1 and v2. (a) to (h) are the noisy mesh; the results of LOM [10], BI [2], GNF [3], CNR [43], GGNF [7] and *NormalNet*; and the ground truth.