# SurfaceNet+: An End-to-end 3D Neural Network for Very Sparse Multi-view Stereopsis

Mengqi Ji*, Jinzhi Zhang*, Qionghai Dai, Lu Fang§

Tsinghua University

**Abstract**—Multi-view stereopsis (MVS) tries to recover the 3D model from 2D images. As the observations become sparser, the significant 3D information loss makes the MVS problem more challenging. Instead of only focusing on densely sampled conditions, we investigate sparse-MVS with large baseline angles since the sparser sensation is more practical and more cost-efficient. By investigating various observation sparsities, we show that the classical depth-fusion pipeline becomes powerless for the case with a larger baseline angle that worsens the photo-consistency check. As another line of the solution, we present SurfaceNet+, a volumetric method to handle the 'incompleteness' and the 'inaccuracy' problems induced by a very sparse MVS setup. Specifically, the former problem is handled by a novel volume-wise view selection approach. It owns superiority in selecting valid views while discarding invalid occluded views by considering the geometric prior. Furthermore, the latter problem is handled via a multi-scale strategy that consequently refines the recovered geometry around the region with the repeating pattern. The experiments demonstrate the tremendous performance gap between SurfaceNet+ and state-of-the-art methods in terms of precision and recall. Under the extreme sparse-MVS settings in two datasets, where existing methods can only return very few points, SurfaceNet+ still works as well as in the dense MVS setting.

**Index Terms**—Multi-view Stereopsis, Volumetric MVS, Sparse Views, Occlusion Aware, View Selection.

---

## 1 INTRODUCTION

MULTI-VIEW stereopsis (MVS) aims to recover a dense 3D model from a set of 2D images with known camera parameters. As the observations become sparser, the more 3D information of the imaged scene get lost during the sensing procedure, making the following perception procedure, for example, an MVS task, more challenging. Dense multi-view sensation has attracted tremendous attention in light field imaging and rendering. Its advantages, such as being robust to occlusion [1] [2] and reducing image noise [3] [4], have been well studied. Unfortunately, it is impractical to densely sample a scene for high-resolution 3D reconstruction, especially for the large-scale scenes. In contrast, the sparser sensation with a wide baseline is more practical and more cost-efficient; however, it aggravates the difficulty of MVS problem since the larger baseline angles lead to tough dense-correspondence matching.

We propose an imperative sparse-MVS leader-board and call for the community's attention on the *general* sparse MVS problem with a large range of baseline angle that could be up to $70°$. Despite of several approaches recovering 3D model from a single view, they are biased towards recovering specific objects or scenes with poor generalization ability. For instance, some work focus on improving the depth map generation with the aid of semantic embeddings [11] [12] [13] or object-level shape prior [14] [15] [16]. Other methods [9] [17] [7] [18] [8] [19], classified as depth
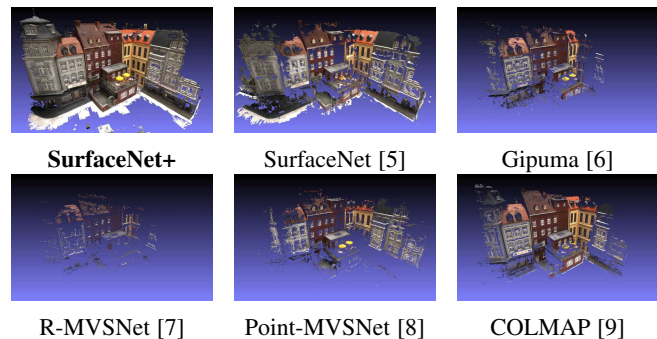


Fig. 1: Illustration of a very sparse MVS setting using only one seventh of the camera views, i.e., $\{v_i\}_{i=1,8,15,22,...}$, to recover the model 23 in the DTU dataset [10]. Compared with the state-of-the-art methods, the proposed SurfaceNet+ provides much complete reconstruction, especially around the boarder region captured by very sparse views.

map fusion algorithms, try to estimate the depth map for each camera view and fuse them into a 3D model. Unfortunately, for the sparse MVS setting with the large baseline angle, e.g. larger than $10°$, these algorithms suffer from incomplete models, because the large baseline angle leads to significantly skewed matching patches from different views and worsens the photo-consistency check. Additionally, as the baseline angle gets larger, the 2D regularization on the depth maps is less helpful for a complete and smooth 3D surface. Because the 2D observation is formed by uneven samples on the 3D surface, the photo consistency agreements can be hardly met by the depth predictions from two views with the large baseline angle, as shown in Fig. 1 and Fig. 3a.

Instead of fusing multiple 2D information into 3D, for the first

* : *Equal Contribution*
§ : *Corresponding Author (fanglu@sz.tsinghua.edu.cn, http://luvision.net)*
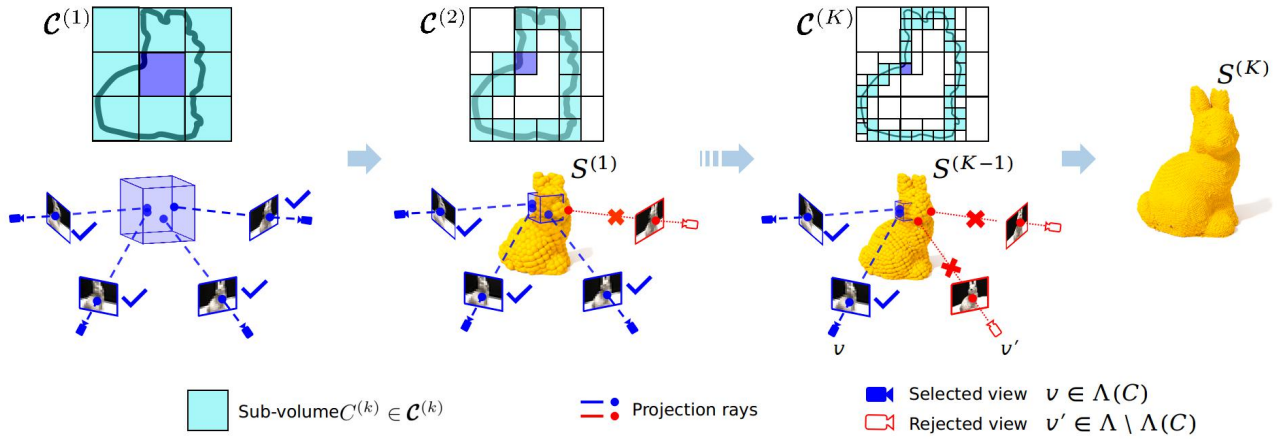
Fig. 2: SurfaceNet+ recovers the whole scene $S^{(K)}$ by progressive refinement of the geometric predictions $S^{(k)}$. So that for each sub-volume $C^{(k)} \in \mathcal{C}^{(k)}$ (drawn as blue cube) the occlusion-aware view selection is performed on the geometric prior. The occluded projection rays are drawn in red and the blue views are the selected ones for reconstruction. In each scale, the volume-wise algorithm only loops through the region in cyan to boost the precision and efficiency.

time, SurfaceNet [5] optimizes the 3D geometry in an end-to-end manner by directly learning the volume-wise geometric context from 3D unprojected color volumes. Even though directly utilizing the 3D regularization may avoid the aforementioned shortcomings of the depth map fusion methods, it still suffers from distinct disadvantages such as noisy surface and large holes around the regions with the repeating pattern and complex geometry. The main reason is that the volume-wise predictions are independently performed without global geometric prior. Consequently, around the region with the repeating pattern, SurfaceNet returns periodic floating surface fragments around the ground-truth surface. Additionally, such noisy predictions further interferes the view selection and leads to large black holes, as shown in Fig. 1.

In this paper, we present an end-to-end learning framework, SurfaceNet+, attacking the very sparse MVS problem. As the sensation sparsity increase, the number of available photo-consistent views becomes less and the view selection scheme gets more critical. Therefore, to adapt to a large range of degree of sparsity, the core innovation is a trainable occlusion-aware view selection scheme that takes the geometric prior into account via a coarse-to-fine scheme. Such volume-wise view selection strategy can significantly boost the performance of the learning-based volumetric MVS methods. More specifically, as shown in Fig. 2, it starts from very coarse 3D surface prediction using all the view candidates, and consequently refines the recovered geometry by gradually discarding the occluded views based on the coarser level geometric prediction. Unlike the traditional image-wise [20] [21] [22] or pixel-wise view selection [9], which cannot filter out the less irrelevant visible views for the final 3D model fusion, the proposed occlusion-aware volume-wise view selection can identify the most valuable view pairs for each 3D sub-volume and the ranking weights is end-to-end trainable. Therefore, consequently only a little proportion of view pairs is needed for volume-wise surface prediction with little performance reduction. That can dramatically reduce the computational complexity by removing redundancy of the multiview sampling. Benefited from the coarse-to-fine fashion, SurfaceNet+ makes the volume-wise occlusion detection more feasible and leads to a high-recall 3D model.

The proposed sparse-MVS leader-board is built on the large-scale DTU dataset [10] and the Tanks-and-Temples dataset [23] with sparsely sampled camera views. The sparse-MVS setting selects one view from every $n$ consecutive camera index, i.e., $\{1, n + 1, 2n + 1, ...\}$, where $n$ is termed as *Sparsity* positively related with baseline angle. The poor performance of the state-of-the-art MVS algorithms on the proposed leader-board demonstrates the necessity of further effort and attention from the community on achieving MVS with various degrees of sparsity. Additionally, the extensive comparison depicts the tremendous performance improvement of SurfaceNet+ over existing methods in terms of precision, recall, and efficiency. As illustrated in Fig. 1, under a very sparse camera setting, SurfaceNet+ predicts a much more complete 3D model compared with recent methods, especially around the border region viewed by a less number of cameras. In summary, the technical contributions in this work are twofold.

- In consideration of the practical necessity of very sparse MVS and the poor performance of the existing MVS methods, we propose a sparse MVS evaluation benchmark and call for the community's attention on the general sparse MVS problem with a broad range of baseline angles.
- To tackle with the sparse MVS problem, we propose a novel trainable occlusion-aware view selection scheme, which is a volume-wise strategy and can significantly boost the performance of the volumetric MVS learning framework. The benchmark and the implementation are publicly available at https://github.com/mjiUST/SurfaceNet-plus.

## 2 RELATED WORK

### 2.1 Multi-view Stereopsis Reconstruction

Works in the multi-view stereopsis (MVS) field can be roughly categorised into 1) direct point cloud reconstructions 2) depth maps fusion algorithms and 3) volumetric methods. Point-cloud-based methods operate directly on 3D points, usually relying on the propagation strategy to gradually densify the reconstruction
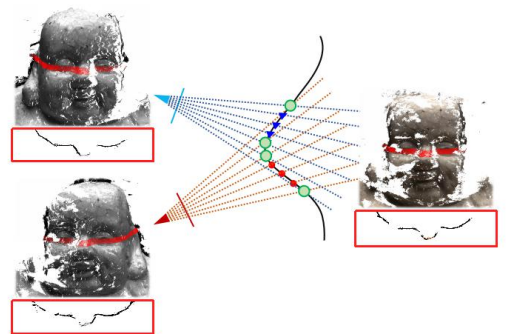
[24] [25]. As the propagation of point clouds proceeds sequentially, these methods are difficult to be fully parallelized and usually take a long time in the processing. Depth maps fusion algorithms [26] [27] [6] decouples the complex MVS problem into relatively small problems of per-view depth map estimation, which focus on only one reference and a few source images at a time and then fuse together with the point cloud [28]. Yet they suffer from incomplete fusion model with large baseline angle or occluded views since skewed patches and uneven samples on the 3D surface in these cases leads to poor quality photo consistency agreements.

Volumetric-based methods divide the 3D space into regular grids and handle the problem in a global coordinate. They use either implicit representation [29] [30] [31] [32] or explicit surface properties [33] [34] [35] [5] [36] [6] to represent and optimize in a global framework. These methods are easy to be parallelized for a multi-view process using a regularization function [30] [29] to minimize errors through all points by gradient descent. Though they are more robust to data noise and outliers, the downside of this representation is the high memory consumption, leading to space discretization error, so they are only applicable to synthetic data with low-resolution inputs [35]. To deal with the small-scale reconstruction problem, these methods either apply the divide-and-conquer strategy [5], or allow a hierarchical multi-scale structure. [32] [37] use an octree representation network to represent both the structure of the octree and the probability of each cell and reconstruct the scene in a coarse-to-fine manner, so that time and space complexities are proportional to the size of the reconstructed model. To perceive more geometry details with limited memory, [38] [39] adopt a hierarchical adaptive multi-scale algorithm and further facilitates the prediction of high-resolution surfaces. Compared with the mentioned volumetric-based methods, the proposed SurfaceNet+ shares the ideal with the divide-and-conquer strategy but infers the 3D surface in a coarse-to-fine fashion with dynamic view-selection strategy.
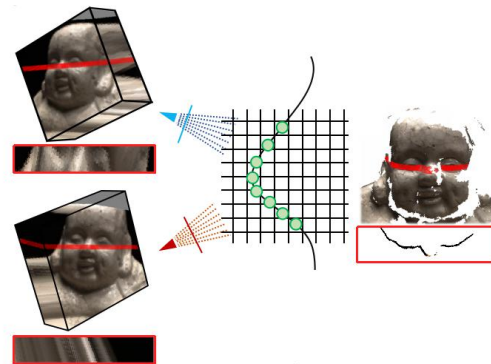
## 2.2 Learning-based MVS

Many learning-based MVS methods have also been developed in recent years. 2D-convolutional neural networks (2D-CNNs) [40] [41] [42] are applied for better patch representation and matching, and others such as [43] learn the normals of a given depth map to improve the depth map fusion. Yet these methods focus on improving the individual steps in the pipeline and their performance is limited in challenging scenes due to the lack of contextual geometry knowledge. The main promotion in this area is 3D cost volume regularization proposed by [44] [18] [35]. This method deploys a 3D volume in the scene space or in the reference camera space. Then, an inverse projection procedure is applied to the 3D volume from several 2D image features gained from different camera positions. Other similar processes such as colored voxel cube [5] and recurrent regularization [7] also use unprojected volumes to get 3D information from 2D image features. The key advantage to process a 3D volume instead of 2D features is that the camera position image information can implicitly write into the 3D volume and the 3D geometry of the scene can be predicted by 3D convolutional layers explicitly. Additionally, during the convolution process, the network is doing the work as in the patch matching method in a highly parallel way, regardless of image distortion and various light conditions. Our approach is more closely related to SurfaceNet [5], which

encodes camera geometries in the network as multiple unprojected volumes to infer the surface prediction in the global coordinate.



(a) Depth map fusion method.



(b) Volumetric method

Fig. 3: Illustration of two types of multi-view reconstruction methods. The front view of the 3D model and the top view of the selected region (red) are shown in pair. The circles (green) indicate the prediction. (a): Because the 2D image unevenly samples the 3D surface, as the baseline angle increases, it is rare for view pair (red and blue) to have intersected rays during depth fusion. The 2D regularization gets less helpful. (b): Volumetric method optimizes the 3D geometry by directly regularizing in 3D space.

## 2.3 Depth Map Fusion Methods

The depth map fusion algorithms first recover depth maps [45] from view pairs by matching similarity patches [20] [21] [22] along the epipolar line and then fuse the depth maps to obtain a 3D reconstruction of the object [26] [27] [6]. [26] is designed for ultra high-resolution image sets and uses a robust descriptor for efficient matching purposes. In [46] describes a patch model that consists of a quasi-dense set of rectangular patches covering the surface. To aggregate image similarity across multiple views, [6] obtains more accurate depth maps. However, in views with the large baseline angle it is problematic with the photo-consistency check because of the significantly skewed patches from different view angles. Therefore, it suffers from incomplete models in sparse-MVS.

After getting multiple depth maps, the depth map fusion algorithm integrates them into a unified and augmented scene representation while mitigating any inconsistencies among individual estimates. To improve fusion accuracy, [27] learns several sources of the depth map outliers. [47] fuses multiple depth estimations into a surface by evaluating visibility in 3D space, and also attempts to reconstruct the region that is not directly

supported by depth measurements. [48] proposes to explicitly target the reconstruction from crowd-sourced images. [29] proposes a variational depth map formulation that enables parallelized computation on the GPU. COLMAP [9] directly maximizes the estimated surface support in the depth maps and allows dataset-wide, pixel-wise sampling for view selection. However, as the observations become sparser, 2D depth fusion regularization is less helpful for a complete 3D model, because each 2D view is formed by uneven samples on the 3D surface and the sparse MVS scenario can hardly lead to photo-consistency agreements of the 3D surface prediction from multiple views. Compared with the heuristic pixel-wise and image-wise view selection methods that manually filter out the occluded views, the proposed volume-wise view selection method is end-to-end trainable from both geometric and photometric priors for each sub-volume.

## 2.4 Review SurfaceNet

SurfaceNet [5] firstly proposes an end-to-end learning framework for MVS by automatically learning both photo-consistency and geometric relations of the surface structure. Given two images $(I_i, I_j)$ and the corresponding camera views $(v_i, v_j)$, SurfaceNet reconstructs the 3D surface in each sub-volume $C$ by estimating for each voxel $x \in C$ whether it is on the surface or not.

Firstly, each image of $I_i$ and $I_j$ is unprojected into $C$ by colorizing the voxels on a traced pixel ray into the same pixel color, so that the new representation $(I_i^C, I_j^C)$ encodes the camera parameters implicitly. The gleaming point of the unprojected sub-volume is view-invariant, because the sub-volume is under the global coordinate rather than the relevant coordinate, like the view-variant sweep plane widely used by depth-fusion methods [8] [7]. So that it does not lead to the uneven sampling effect.

Then, a pair of colored voxel cubes $(I_i^C, I_j^C)$ is fed into SurfaceNet, a fully 3D convolutional neural network, to predict for each voxel $x \in C$ the confidence $p_x \in (0, 1)$, which indicates whether a voxel is on the surface or not by using cross-entropy loss. Due to the fully convolutional design, the sub-volume size $s^3$ for inference can be different from that for training, and can be adaptive to various graphic memory sizes.

Lastly, to generalize to a case with multiple views $\Lambda = \{v_1, ..., v_i, ..., v_j, ..., v_V\}$, it only selects a subset of view pairs $(v_i, v_j)$ to predict $p_x^{(v_i, v_j)}$, i.e., the confidence that a voxel $x$ is on the surface, then combines together by taking the weighted average of the predictions based on the relative weight $w_C^{(v_i, v_j)}$ for each view pair

$$w_C^{(v_i, v_j)} = r\left(\theta_C^{(v_i, v_j)}, e(C, I_{v_i}), e(C, I_{v_j})\right), \quad (1)$$

which is inferred by function $r(\cdot)$ with the inputs of the patch embeddings $e(\cdot)$ and the baseline angle $\theta_C^{(v_i, v_j)}$, i.e., the angle between the projection rays from the center of $C$ to the optical centers of $v_i$ and $v_j$. So that the volume-wise reconstruction becomes computationally feasible by ignoring the majority of possible view pairs.

Benefited from the direct regularization of the 3D surface, SurfaceNet does not suffer from the shortcoming of 2D regularization owing to the uneven sample of 2D projection. However, the view selection scheme becomes non-trivial and is challenging for the sparse MVS scenario where SurfaceNet still has distinct disadvantages, such as large holes and noisy surfaces around the regions with complex geometry and repeating patterns. Additionally, the volume-wise prediction becomes extremely computationally heavy for large scene reconstruction. In this paper, SurfaceNet+ solves the aforementioned problems with a large margin of performance improvement and around 10X speedup compared with SurfaceNet.

## 3 SURFACENET+

In this Section, We present SurfaceNet+, an end-to-end learning framework, to handle the very sparse MVS problem, where the critical problem to be solved is the view selection. As the sensation sparsity increases, the number of available photo-consistent views becomes less; thus, the view selection scheme gets more critical. SurfaceNet+ utilizes a novel trainable occlusion-aware view selection scheme that takes the geometric prior into account via a coarse-to-fine strategy. In short, the multi-scale inference (subsection 3.1) outputs the geometric prior required by the occlusion-aware view selection scheme (subsection 3.2). As shown in Fig. 2, starting from a bounding box, a very coarse 3D surface is predicted by considering all the view candidates. Subsequently, the coarse level geometry gets iteratively refined by gradually discarding the occluded views based on the coarser level geometric prior. In subsection 3.3, the backbone network, a fully convolutional network structure, is presented in detail.

### 3.1 Multi-scale Inference

For a volume-wise reconstruction pipeline, the noisy prediction occurs frequently around the 3D surface with repeating patterns. Moreover, it suffers from a huge computational burden to iterate through the majority of the empty space. While it may be intuitive to consider the 3D geometry prior during reconstruction, the difficulty lies in that the general MVS task does not have any shape prior of the scene. What we propose is a coarse-to-fine architecture to gradually refine the geometric details under the assumption that the minority volume of the space is occupied by the 3D surface of the scene.

In the first stage, SurfaceNet+ divides the entire bounding box into a set of sub-volumes $\boldsymbol{C}^{(1)}$ of the coarsest level with the side length $l^{(1)} = s \cdot r^{(1)}$, where $r^{(1)}$ is the voxel resolution of the coarsest level when the voxelization forms a tensor of size $s \times s \times s$. The tensor size depends highly on the graphic memory size, for example $s = 32/64/128$. As the output of this stage, the estimated surface of the coarsest level is denoted as $S^{(1)}$, where $x \in S$ means an occupied voxel in the surface prediction.

The following iterative stage divides the space into different sub-volume set in each scale level, i.e., $\{\boldsymbol{C}^{(2)}, \cdots, \boldsymbol{C}^{(k)}, \cdots, \boldsymbol{C}^{(K)}\}$, whose resolutions are a geometric sequence with the common ratio $\delta$, i.e., $r^{(k)} = \delta \cdot r^{(k+1)}$. Usually, we set $\delta = 4$ to compromise between efficiency and effectiveness. This procedure is iterated until meeting the condition $r^{(K)} \leq r$, where $r$ is the desired resolution and $r^{(K)}$ is the finest one. The way to divide the sub-volume is highly dependent on the predicted point cloud of the coarser level $S^{(k-1)}$, when $k = 2, 3, \cdots$, so that each of the regular sub-volume divisions $C^{(k)} \in \boldsymbol{C}^{(k)}$ contains at least one point:

$$\boldsymbol{C}^{(k)} = \arg\min_{\boldsymbol{C}}\{|\boldsymbol{C}| \,|\forall \boldsymbol{C} : \quad (2)$$
$$(S^{(k-1)} \subseteq \bigcup \boldsymbol{C}) \wedge (\forall C \in \boldsymbol{C} : S^{(k-1)} \cap C \neq \varnothing)\},$$
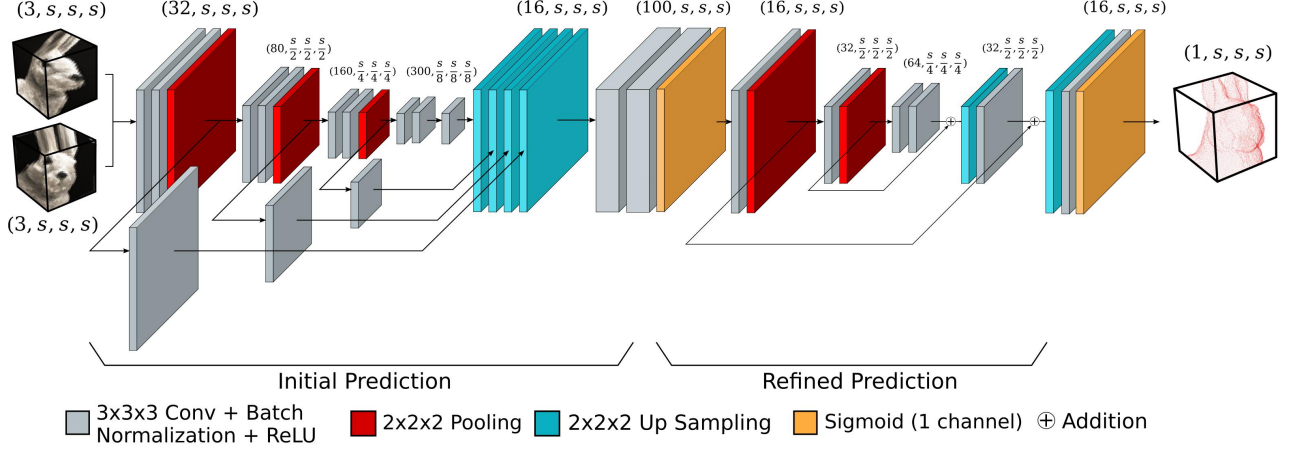
Fig. 4: The network design of SurfaceNet+. The input of the network is two unprojected sub-volumes with size of (3,s,s,s) from different views. The final prediction is an one channel tensor predicting for each voxel the probability of being on surface.

where $|\mathcal{C}|$ denotes the number of sub-volume divisions, and $\bigcup \mathcal{C}$ is a short representation for the union of all the sub-volumes, i.e., $\bigcup_{C \in \mathcal{C}} C$. To eliminate the boundary effect of the convolution operation, we usually loose the above limitation and allow a slight overlapping between the neighboring sub-volume. The point cloud output $S^{(k)}$ of SurfaceNet+ will be introduced in subsection 3.3.

### 3.2 Trainable Occlusion-aware View Selection

As depicted in Fig. 1, even though SurfaceNet [5] does not have the artifacts caused by uneven sampling from 3D surface to 2D depth, it suffers from large holes around the complex geometry. The key reason is that the view selection becomes more critical for the sparse MVS problem. Following the annotation in subsection 2.4, we introduce how the proposed trainable occlusion-aware view selection scheme can rank and select the top-$N_v$ most valuable view pairs $\boldsymbol{V}_C$ of each sub-volume $C$ from all the possible view pairs

$$\boldsymbol{V} = \{(v_i, v_j) | (v_i, v_j \in \Lambda) \wedge (v_i \neq v_j)\}, \quad (3)$$

based on the learned relative weights $w_C^{(v_i, v_j)}$, which is inferred from both the geometric and photometric priors. Note that the multi-scale scheme can provide us with the crucial geometric prior $S^{(k-1)}$. Consequently, according to Eq. 8, the surface in each sub-volume $C$ is fused by the $|\boldsymbol{V}_C| = N_v$ predictions.

**Geometric Prior.** The geometric prior can be easily encoded from the multi-scale predictions. For any camera view $v$ w.r.t. each sub-volume $C \in \mathcal{C}$, a convex hull $H(C, v) \subset \mathbb{R}^3$ is uniquely defined by a set of points

$$H(C, v) = Conv(\Gamma(C) \cup \{o_v\}), \quad (4)$$

where $o_v$ is the camera center of $v$, and the set $\Gamma(C) = \{c_1, c_2, ..., c_8\}$ contains the 8 corners of $C$.

The more points in the coarser level of surface prediction $S^{(k-1)}$ that appear in the region between the camera view $v$ and the sub-volume $C^{(k)}$, the more likely the view $v$ is occluded. These barrier points are defined in the set

$$B(C^{(k)}, v) = S^{(k-1)} \cap H(C^{(k)}, v) \backslash C^{(k)}. \quad (5)$$

**Trainable Relative Weights.** As suggested in [5], the end-to-end trainable relative weights not only can improve the efficiency

by filtering out the majority of the less crucial view pairs for each sub-volume but also can improve the effectiveness of the surface prediction by weighted fusion. Note that, for sparse MVS, the number of the valid views for each sub-volume could be too few to heuristically detect occlusions. Instead, we propose a trainable occlusion-aware view pair selection scheme that learns the relative weights based on both the geometric and photometric priors:

$$w_C^{(v_i, v_j)} = p_{C^{(k)}}^{(v_i, v_j)} \cdot r\left(\theta_C^{(v_i, v_j)}, e(C, I_i), e(C, I_j)\right), \quad (6)$$

where, the photometric priors are the same as SurfaceNet [5], i.e., the baseline angle $\theta_C^{(v_i, v_j)} = \angle o_{v_i} o_C o_{v_j}$ as well as the embeddings $e(\cdot)$ of the cropped patches around the 2D image of $o_C$ on both $I_i$ and $I_j$ in Eq. 1, and the geometric prior is encoded as the probability of being not occluded, i.e., :

$$p_{C^{(k)}}^{(v_i, v_j)} = \exp\left(-\alpha \cdot r_k^2 \cdot (|B(C^{(k)}, v_i)| + |B(C^{(k)}, v_j)|)\right), \quad (7)$$

where $\alpha$ is a hyper parameter controlling the sensitivity of this occlusion probability term and the coefficient $r_k^2$ can be understood as a normalization term w.r.t. different scales. In Section we will show the effect of $\alpha$ and how it improves the performance of the reconstruction.

**Weighted Average Surface Prediction.** Lastly, for the general MVS problem, we follow the fusion strategy in SurfaceNet [5], which ranks and selects only a small subset of view pairs $\boldsymbol{V}_C$. Subsequently, the confidence that a voxel $x$ is on the surface, $p_x$, is inferred by the weighted average of the predictions $p_x^{(v_i, v_j)}$:

$$p_x = \frac{\sum_{(v_i, v_j) \in \boldsymbol{V}_C} w_C^{(v_i, v_j)} p_x^{(v_i, v_j)}}{\sum_{(v_i, v_j) \in \boldsymbol{V}_C} w_C^{(v_i, v_j)}}, \quad (8)$$

where $\boldsymbol{V}_C$ denotes the set of selected view pairs with the size of $|\boldsymbol{V}_C| = N_v$, and the relative weight $w_C^{(v_i, v_j)}$ for each view pair is end-to-end trainable and is inferred by Eq. 6. Note that a smaller $N_v$ can lead to more efficient and less effective results, which is discussed in section 5.

### 3.3 Network

**Network Architecture.** At each stage of reconstruction, we use a 3D convolutional neural network to predict whether each voxel

in each sub-volume is on the surface or not. Specifically, given $C^k$ and the corresponding image view pairs $(I_i, I_j)$, we first blur each image using a Gaussian kernel to spread the local information around the large receptive field and to guarantee the image consistency in all stages. The unprojected 3D sub-volume $(I_i^{C^{(k)}}, I_j^{C^{(k)}})$ for a view pair is demonstrated in Fig. 3b. The beauty of this representation is that it implicitly encodes the camera parameters as well as scale information to adapt to a fully convolutional neural network.

The detailed network configuration is shown in Fig. 4. The building blocks of the model are a UNet-like architecture followed by a refinement network. SurfaceNet+ takes two colored sub-volumes $(I_i^{C^{(k)}}, I_j^{C^{(k)}})$ as input, which stores two RGB color values and forms a 6-channel tensor of shape $6 \times s \times s \times s$, and predicts the on-surface probability for each voxel $p_{x \in C^{(k)}}$ forming a tensor of size $1 \times s \times s \times s$, To extract distinct geometry information in various scales, we first use a pyramid structure to process the features in different receptive fields. To better aggregate multi-scale information, we use two $3 \times 3$ convolution layers followed by a one-channel convolution layer with a sigmoid activation function after concatenating the features on different scales. Inspired by [18], we apply a prediction refinement network at the end of the previous network. After the initial output $\tilde{S}_C^{(k)}$ with a tensor shape of $1 \times s \times s \times s$, the skip connections at each layer are used to learn the residual prediction and to generate the final output $S_C^{(k)}$.

**Loss.** The training loss consists of two parts to penalize both the initial prediction $p_x$ and the refined prediction $p'_x$. In the first stage, the discriminative prediction per voxel $p_x$ is compared with the ground-truth $\hat{s}_x$. Since the majority of the voxels does not contain the surface, i.e., $\hat{s}_{x \in C^{(k)}} = 0$, a class-balanced cross-entropy function is utilized, i.e., for each $C^{(k)}$ we have

$$L_{init} = \qquad\qquad\qquad\qquad\qquad\qquad (9)$$
$$- \sum_{x \in C^{(k)}} \left\{ \beta^{(k)} \hat{s}_x \log p_x + (1 - \beta^{(k)})(1 - \hat{s}_x) \log(1 - p_x) \right\},$$

where the hyper-parameter $\beta^{(k)}$ is the occupancy ratio of the ground-truth in the scale $k$.

In the second stage, the refined prediction $p'_x$ is regressed to the ground-truth by the mean square error (MSE), so that the small residue can be penalized as well,

$$L_{refine} = \sum_{x \in C} \|\hat{s}_x - p'_x\|_2, \qquad (10)$$

where $p_x \in S_C^{(k)}$. Consequently, the training loss is defined as:

$$L_{total} = L_{refine} + L_{init}. \qquad (11)$$

### 3.4 Implementation Details

Our network is trained on the DTU dataset [50]. We use the volume with $32^3$ voxels to train the network, with a batch size of 16, and the voxel resolution is separately set to 0.4*mm*, 1.0*mm* and 2.0*mm* for each set to generalize on a different scale of surface geometry. To acquire a favorable generalization on sparse-MVS, the network needs to be trained from a variety of view pairs. Therefore, the 3D convolutional network is first trained on randomly-sampled non-occluded view pairs $(v_i, v_j)$ without relative weight $w_C^{(v_i, v_j)}$. Then the training process is combined together with $w_C^{(v_i, v_j)}$, and the view pair number is

fixed to 6. Specifically, the relative weights learning procedure is performed using a 2-layer fully connected neural network $r(\cdot)$. The computation is introduced in subsection 3.2 except that the surface prediction at the previous stage is replaced with the reference model. During the reconstruction stage, the volume size is $64^3$ and the output is upsampled to $128^3$. All the training and reconstruction processes are accomplished on one GTX 1080Ti graphics card.

## 4 SPARSE-MVS BENCHMARK

In this section, the imperative sparse-MVS leader-board on different datasets, the DTU dataset [10], the Tanks and Temples dataset [23] (T&T), and the ETH3D low-res dataset [49], is introduced with extensive comparisons to the recent MVS methods under various observation sparsity levels.

We benchmark SurfaceNet+ at all sparsities from 1 to 11 against several state-of-the-art methods. The sparse MVS setting in our leader-board selects a small proportion of the camera views by consecutively sampling a view $v_i$ from every $sparsity = n$ camera index, i.e., $\{1, n + 1, 2n + 1, \cdots\}$. In reality, it is also practical to sample small-batches of images at sparse viewpoints, i.e., grouping batches of views with certain Batchsize at the previously defined sparse viewpoints with a certain Sparsity. When $Sparsity = 3$ and $Batchsize = 1$, the chosen camera indexes are 1 / 4 / 7 / 10 / $\cdots$. When $Sparsity = 3$ and $Batchsize = 2$, the chosen camera indexes are 1,2 / 4,5 / 7,8 / 10,11 / $\cdots$.

Fig. 5 depicts the relationship between sparsity $n$ and the average baseline angle $\bar{\theta}$ averaging over all the ground-truth points in the 22 models of the DTU dataset, 8 models of the Tanks and Temples dataset, and 5 models of the ETH3D low-res dataset, respectively. Note that, for simplicity, only the nearest view pairs are considered to calculate the baseline angle statistics.

$$\theta = \{ \angle o_{v_i} x o_{v_j} | x \in \hat{S}, v_i \in \Lambda, v_j = \arg \min_{v \in \Lambda} \overline{o_{v_i} o_v} \} \qquad (12)$$

As the sparsity increases $n = 1, ..., 11$, the average baseline angle $\bar{\theta}$, defined by the intersected projection rays, gradually grows in a large range, e.g. reaching more than $70°$ in both DTU and T&T datasets. Due to the positive correlation between $n$ and $\bar{\theta}$, we claim that our sparse-MVS setting is reasonable by not only covering various degrees of sparsity but also containing irregular sampling locations.

### 4.1 DTU Dataset [10]

We qualify the performances on the DTU dataset [10] in different sparse MVS settings. The DTU dataset is a large-scale MVS benchmark, which features a variety of objects and materials, and contains 80 different scenes seen from 49 camera positions under seven different lighting conditions. 22 models are selected from the DTU dataset as the evaluation set, following [5] [1].

The chart in Fig. 6 plots the performance under a large range of sparsity in terms of *f-score* (1*mm*), which unifies both *recall* and *precision*. This apparently shows that our proposed method

---

1. Follow the same dataset split in SurfaceNet [5]. Training: 2, 6, 7, 8, 14, 16, 18, 19, 20, 22, 30, 31, 36, 39, 41, 42, 44, 45, 46, 47, 50, 51, 52, 53, 55, 57, 58, 60, 61, 63, 64, 65, 68, 69, 70, 71, 72, 74, 76, 83, 84, 85, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 107, 108, 109, 111, 112, 113, 115, 116, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128. Validation: 3, 5, 17, 21, 28, 35, 37, 38, 40, 43, 56, 59, 66, 67, 82, 86, 106, 117. Evaluation: 1, 4, 9, 10, 11, 12, 13, 15, 23, 24, 29, 32, 33, 34, 48, 49, 62, 75, 77, 110, 114, 118
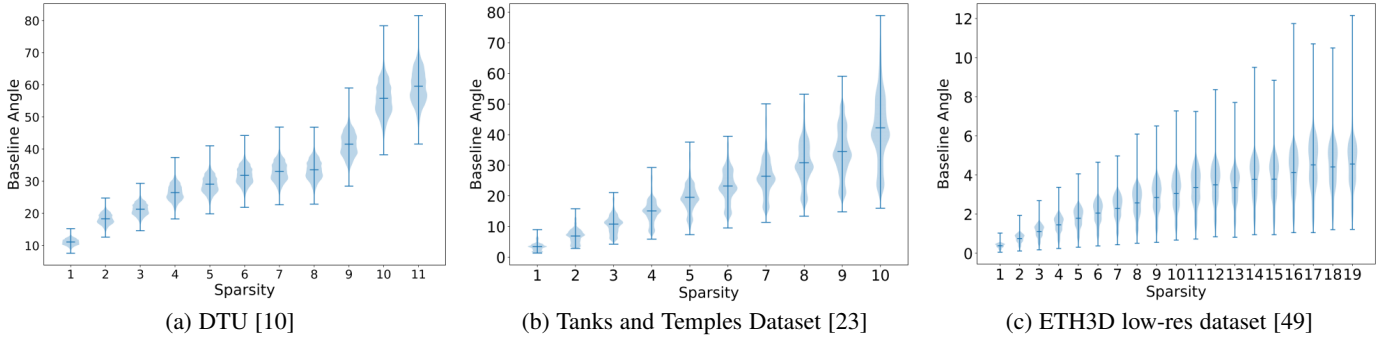
Fig. 5: The relationship between sparsity and the average baseline angle over all the models in the DTU dataset [10], the Tanks and Temples dataset [23] and the ETH3D low-res dataset [49].
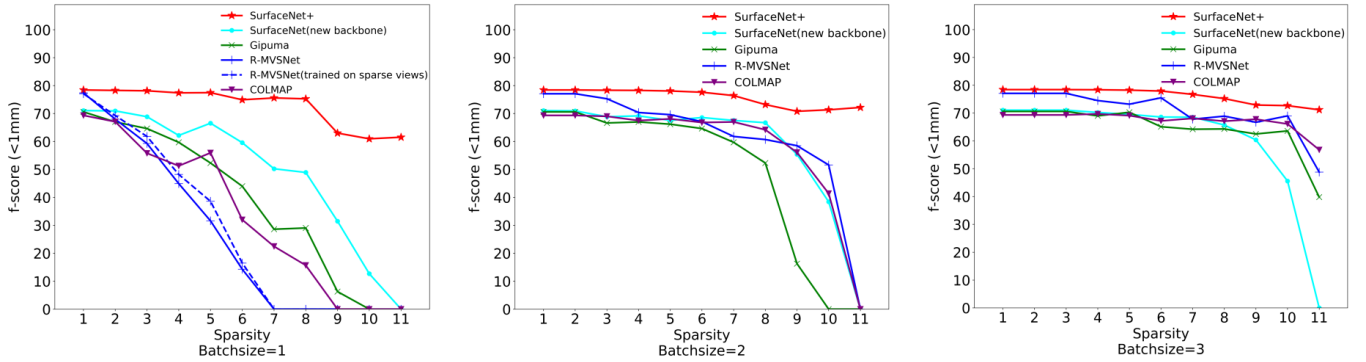


Fig. 6: Comparison with the existing methods in the DTU Dataset [10] with different sparsely sampling strategy. When $Sparsity = 3$ and $Batchsize = 2$, the chosen camera indexes are 1,2 / 4,5 / 7,8 / 10,11 / .... SurfaceNet+ constantly outperforms the state-of-the-art methods at all the settings, especially at the very sparse scenario.

constantly outperforms others in all the sparse settings. Especially for the case of $\bar{\theta} < 40°$, amazingly, SurfaceNet+ constantly performs well without obvious degradation. In the extremely sparse case, i.e., $\bar{\theta} > 50°$, as expected, SurfaceNet+ shows a tiny performance reduction. In contrast, other methods, especially the depth-fusion methods, merely predict a few points. Readers can refer to subsection 2.3 for the discussion why the depth-fusion methods cannot return a complete result. In our leader-board, depth-map-based methods such as R-MVSNet [7] and Gipuma [6] share the same depth fusion code. For fair comparison, we tuned the hyper-parameters in the depth fusion algorithm to induce better performance in terms of f-score under $1mm$ at each sparsity setting.

More detailed quantitative results are listed in Table 1, where 3 different matrices are adopted for evaluation. The precision and recall have two metrics: the distance metric [10] and the percentage metric [23]. The *overall* score for the percentage metric is measured as the f-score, and a similar measurement for the distance metric overall is given by the average of the mean precision and mean recall. Obviously, SurfaceNet+ outperforms the state-of-the-art methods in both recall and precision at all sparsity settings. Unlike other methods whose recall dramatically decay when the sparsity increases, SurfaceNet+ has almost consistent recall quality with high precision.

The qualitative comparison of SurfaceNet+ and the other two methods, R-MVSNet [7] and Gipuma [6], is illustrated in Fig. 7, showing that SurfaceNet+ precisely reconstructs the scenes while maintaining high recall. In particular, SurfaceNet+ is able to generate a high-recall point cloud in complex geometry and repeating pattern regions when $sparsity = 7$, which means it evenly fuses the accurate 3D model with corrected-selected non-occluded views. The detailed analysis is shown in Section 5.

To have a slightly different way of sparse sampling, three $Batchsize$ values $\{1, 2, 3\}$ are evaluated as depicted in Fig. 6. It can be observed that SurfaceNet+ constantly outperforms others despite that the depth-fusion methods (Gipuma [6], R-MVSNet [7], COLMAP [9]) boost the performance as the $Batchsize$ increases. Moreover, as the disparity increases, the performance drop of the existing methods is apparently larger than that of SurfaceNet+. In particular, we have retrained the R-MVSNet for sparse MVS with randomly-sampled non-occluded view pairs at $Batchsize = 1$. As shown in Fig. 6, the gain is inapparent in terms of f-score. As the depth-fusion based MVS methods (R-MVSNet) rely more on the photo-consistency in 2D images, the large baseline angles of a very sparse MVS problem leads to severely skewed matching patches across views that significantly toughen the dense correspondence problem. In contrast, the learning-based volumetric MVS methods like SurfaceNet+ avoids the 2D correspondence search problem by directly inferring 3D surface from each unprojected 3D sub-volumes. That may explain why the learning-based volumetric methods outperform the depth-fusion based methods in the very sparse MVS settings. For the experiment settings, both R-MVSNet and Gipuma shared the same depth fusion code, and we tuned the hyper-parameters of it to induce better performance in terms of f-score under 1mm at each sparsity setting. More specifically, followed by Gipuma [6], since

Fig. 7: Quanlitative results of three scans 1, 23 and 114 of the DTU dataset compared with R-MVSNet [7] and Gipuma [6]. SurfaceNet+ shows superior performance, particularly with its stable recall quality in sparse cases. Note that the reconstruction of SurfaceNet+ corresponds to the highest completeness and overall quality as seen in Fig. 6 and Table. 1.

Fig. 8: Results of three models in Tanks and Temples 'intermediate' set [23] compared with R-MVSNet [7] and COLMAP [9], which demonstrate the power of SurfaceNet+ of high recall prediction in sparse-MVS.

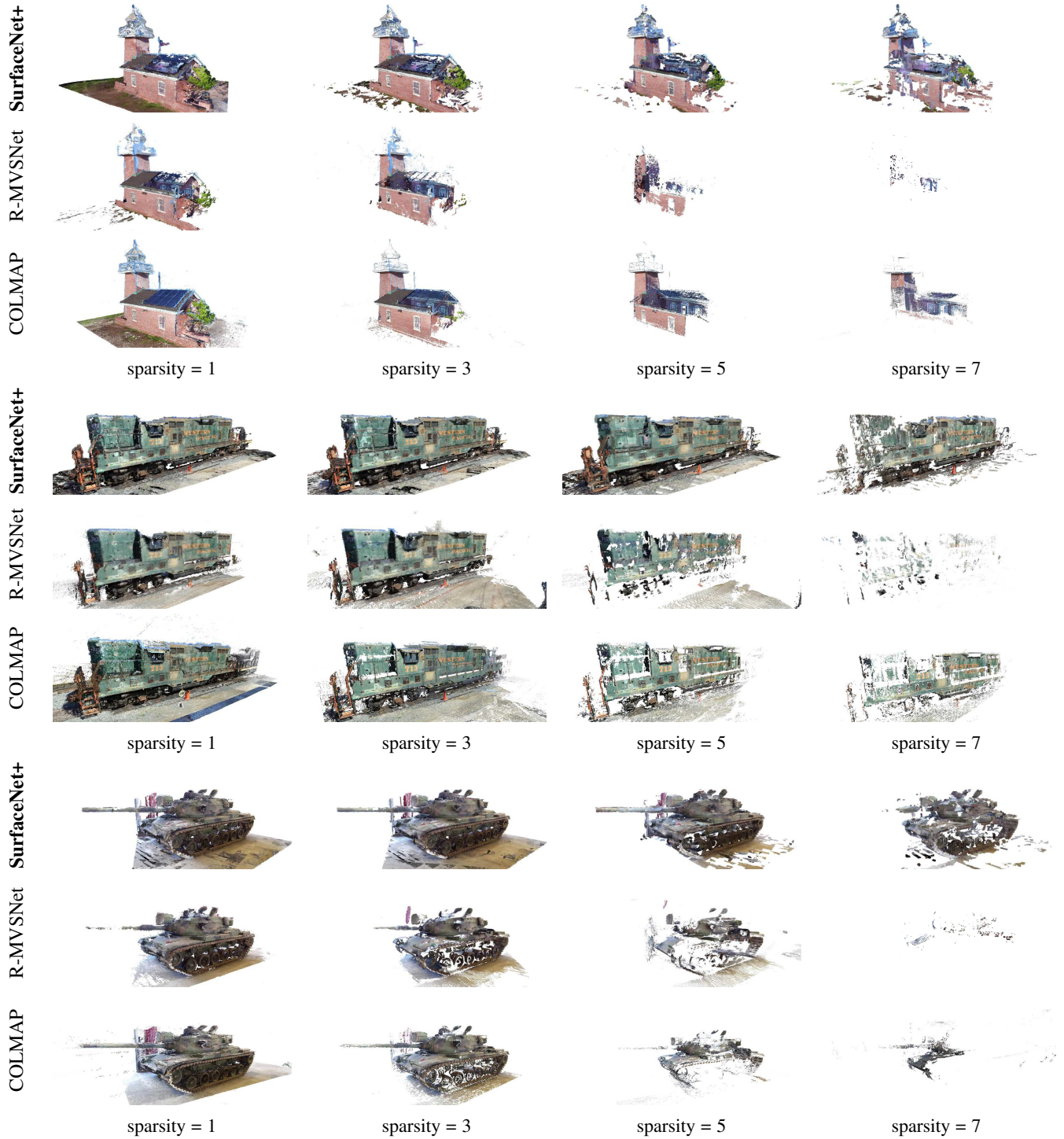| Sparsity | Method | Mean Distance(mm) | | | Percentage(<1mm) | | | Percentage(<2mm) | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Precision | Recall | **Overall** | Precision | Recall | **f-score** | Precision | Recall | **f-score** |
| 1 | SurfaceNet+ | 0.385 | **0.448** | **0.416** | 88.01 | **73.01** | **78.44** | 92.33 | **78.1** | **83.55** |
| | SurfaceNet [5] | 0.450 | 1.021 | 0.735 | 84.49 | 64.58 | 71.65 | 89.10 | 68.72 | 76.21 |
| | Gipuma [6] | **0.283** | 0.873 | 0.578 | **94.65** | 59.93 | 70.64 | **96.42** | 63.81 | 74.16 |
| | R-MVSNet [7] | 0.383 | 0.452 | 0.417 | 87.63 | 72.48 | 77.09 | 91.74 | 76.39 | 82.01 |
| | COLMAP [9] | 0.411 | 0.657 | 0.534 | 82.24 | 52.48 | 61.34 | 88.26 | 62.20 | 72.93 |
| 3 | SurfaceNet+ | 0.446 | **0.482** | **0.464** | 86.06 | **74.41** | **78.15** | 90.87 | **78.25** | **82.91** |
| | SurfaceNet | 0.461 | 0.997 | 0.729 | 83.02 | 61.09 | 68.87 | 88.31 | 66.39 | 74.41 |
| | Gipuma | **0.267** | 1.252 | 0.759 | **95.51** | 50.88 | 64.63 | **97.49** | 50.33 | 63.68 |
| | R-MVSNet | 0.465 | 1.012 | 0.738 | 89.55 | 48.03 | 59.28 | 96.96 | 57.92 | 69.04 |
| | COLMAP | 0.467 | 1.090 | 0.778 | 78.45 | 49.26 | 59.62 | 91.44 | 55.98 | 65.77 |
| 5 | SurfaceNet+ | 0.446 | **0.491** | **0.469** | 88.58 | **71.63** | **77.48** | 92.86 | **76.04** | **82.28** |
| | SurfaceNet | 0.445 | 0.948 | 0.701 | 81.07 | 58.62 | 66.55 | 85.40 | 62.76 | 70.97 |
| | Gipuma | 0.460 | 1.633 | 1.046 | **92.38** | 38.53 | 52.36 | **95.10** | 48.15 | 61.78 |
| | R-MVSNet | **0.329** | 2.209 | 1.269 | 89.26 | 20.51 | 31.60 | 93.99 | 32.74 | 46.37 |
| | COLMAP | 0.443 | 1.284 | 0.863 | 88.79 | 42.51 | 55.94 | 92.91 | 54.89 | 65.77 |
| 7 | SurfaceNet+ | **0.435** | **0.524** | **0.479** | **91.36** | **72.23** | **75.59** | **95.21** | **76.54** | **81.86** |
| | SurfaceNet | 0.688 | 1.130 | 0.909 | 66.86 | 36.91 | 50.24 | 69.21 | 46.91 | 61.70 |
| | Gipuma | 0.569 | 1.770 | 1.169 | 85.35 | 17.91 | 28.66 | 90.78 | 28.00 | 41.31 |
| | R-MVSNet | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| | COLMAP | 0.545 | 1.756 | 1.150 | 59.28 | 15.14 | 22.46 | 80.92 | 31.56 | 41.89 |
| 9 | SurfaceNet+ | **0.441** | **0.895** | **0.668** | **85.99** | **53.16** | **63.01** | **89.86** | **57.63** | **67.86** |
| | SurfaceNet | 1.112 | 2.176 | 1.644 | 35.84 | 29.53 | 31.47 | 38.36 | 34.01 | 35.49 |
| | Gipuma | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| | R-MVSNet | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| | COLMAP | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| 11 | SurfaceNet+ | **0.445** | **0.880** | **0.663** | **85.81** | **51.52** | **61.54** | **90.05** | **55.41** | **65.99** |
| | SurfaceNet | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| | Gipuma | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| | R-MVSNet | empty | empty | empty | empty | empty | empty | empty | empty | empty |
| | COLMAP | empty | empty | empty | empty | empty | empty | empty | empty | empty |

TABLE 1: Quantitative results of reconstruction quality on the DTU dataset in terms of the distance metric(lower is better) and the percentage metric [23](higher is better) with 1$mm$ and 2$mm$ as thresholds. SurfaceNet+ constantly outperforms the state-of-the-arts in all the sparse-MVS settings with $n = 1, 3, 5, 7, 9, 11$.

| Method | **Average Rank** | Mean | Family | Francis | Horse | Lighthouse | M60 | Panther | Playground | Train |
|---|---|---|---|---|---|---|---|---|---|---|
| ACMM [51] | **14.00** | 57.27 | 69.24 | 51.45 | 46.97 | 63.20 | 55.07 | 57.64 | 60.08 | 54.48 |
| CasMVSNet [19] | **15.75** | 56.84 | 76.37 | 58.45 | 46.26 | 55.81 | 56.11 | 54.06 | 58.18 | 49.51 |
| ACMH [51] | **22.25** | 54.82 | 69.99 | 49.45 | 45.12 | 59.04 | 52.64 | 52.37 | 58.34 | 51.61 |
| UCSNet [52] | **22.62** | 54.83 | 76.09 | 53.16 | 43.03 | 54.00 | 55.60 | 51.49 | 57.38 | 47.89 |
| PLC [53] | **24.38** | 54.56 | 70.09 | 50.30 | 41.94 | 58.86 | 49.19 | 55.53 | 56.41 | 54.13 |
| **SurfaceNet+** | **36.12** | 49.38 | 62.38 | 32.35 | 29.35 | 62.86 | 54.77 | 54.14 | 56.13 | 43.10 |
| Dense R-MVSNet [7] | **41.00** | 50.55 | 73.01 | 54.46 | 43.42 | 43.88 | 46.80 | 46.69 | 50.87 | 45.25 |
| VisibilityAwarePointMVSNet [54] | **43.88** | 48.70 | 61.95 | 43.73 | 34.45 | 50.01 | 52.67 | 49.71 | 52.29 | 44.75 |
| Point-MVSNet [55] | **44.38** | 48.27 | 61.79 | 41.15 | 34.20 | 50.79 | 51.97 | 50.85 | 52.38 | 43.06 |
| R-MVSNet [7] | **46.88** | 48.40 | 69.96 | 46.65 | 32.59 | 42.95 | 51.88 | 48.80 | 52.00 | 42.38 |
| MVSNet [18] | **57.50** | 43.48 | 55.99 | 28.55 | 25.07 | 50.79 | 53.96 | 50.86 | 47.90 | 34.69 |
| COLMAP [9] | **60.50** | 42.14 | 50.41 | 22.25 | 25.63 | 56.43 | 44.83 | 46.97 | 48.53 | 42.04 |

TABLE 2: The top and non-anonymous methods on the Tanks and Temples (T&T) dataset [23] leaderboard. The average rank of SurfaceNet+ is higher than that of R-MVSNet [7], MVSNet [18], COLMAP [9], and Point-MVSNet [55].

there is a tradeoff between accuracy and completeness, we choose the depth fusion parameter settings that achieve high accuracy at sparsity=1,2 and high completeness at sparsity>=3. The other part remain the same as the paper of R-MVSNet [7] and Gipuma [6]. In COLMAP [9], all parameters were set as the default values.

dense MVS condition, the overall rank of SurfaceNet+ is still higher than that of R-MVSNet [7], MVSNet [18], COLMAP [9], and Point-MVSNet [55]. Note that we list and compare with all the top and non-anonymous methods on the leaderboard in the following table.

## 4.2 Tanks and Temples Dataset [23]

The Tanks and Temples (T&T) dataset [23] contains real-world large scenes under complex lighting conditions. In Fig. 8, we compare the qualitative results in the Tanks and Temples dataset [23] with R-MVSNet [7] and COLMAP [9]. The results indicate the effectiveness of our proposed method at different sparsities. We submitted and evaluated the SurfaceNet+ results ($Sparsity = 1$) to the online leader-board. As depicted in Table 2, despite the

## 4.3 Generalization on the ETH3D Dataset [49]

We also evaluate the generalization ability by adopting the ETH3D dataset [49], i.e., we direct evaluate the proposed method that trained only on the DTU training dataset without fine-tuning the network. The results of the low-resolution scenes are shown in Fig. 9. It is worth noting that the baseline angle in the ETH3D dataset is tiny among all the camera views because the images were acquired by just rotating the camera with little camera trans-
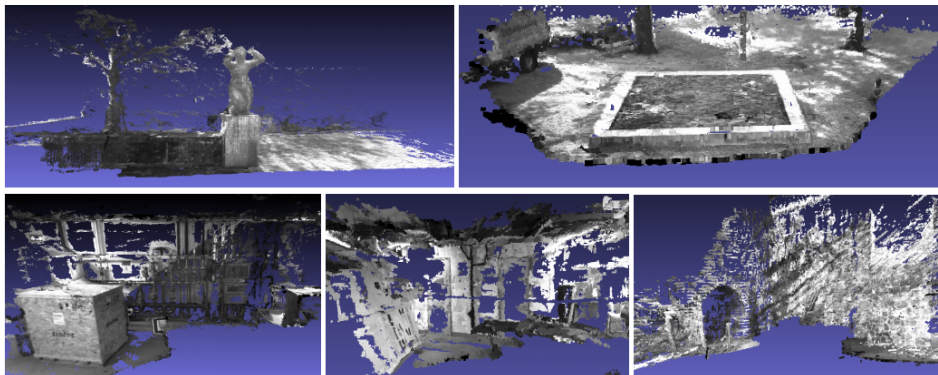
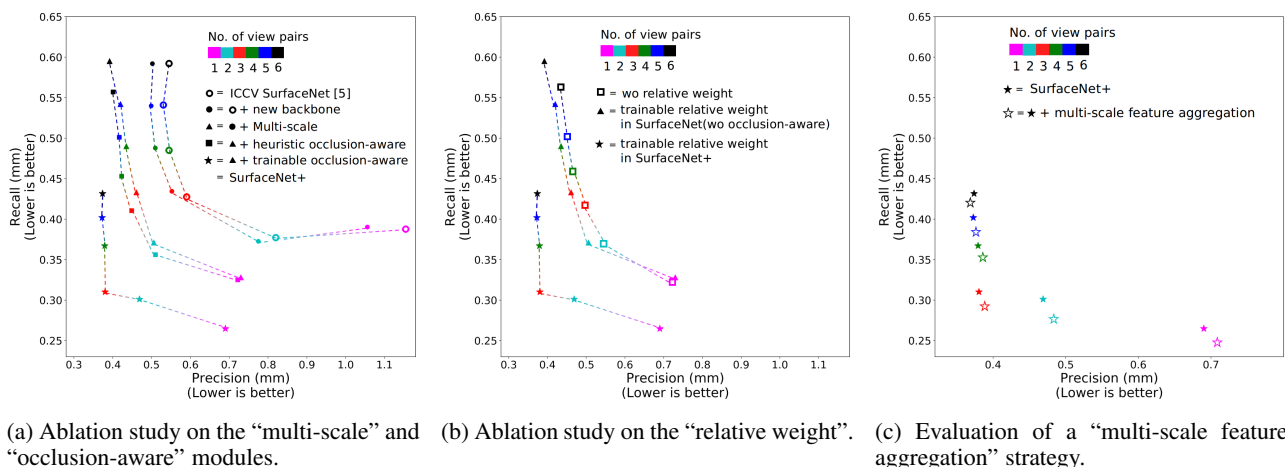Fig. 9: Point cloud reconstructions of the ETH3D low-res dataset [49].



(a) Ablation study on the "multi-scale" and "occlusion-aware" modules.

(b) Ablation study on the "relative weight".

(c) Evaluation of a "multi-scale feature aggregation" strategy.

Fig. 10: Ablation study. (a): Comparison among ICCV SurfaceNet [5] (○ curve), SurfaceNet with new backbone (● curve), **Multi-scale** (▲ curve), heuristic occlusion-aware view selection during inference (■ curve) and the proposed trainable occlusion-aware view selection (★ curve). (b): The benefit from an explicit "relative weight". □ curve indicates the setting without relative weight that takes heuristic occlusion-aware view selection; ▲ curve is the experiment using trainable relative weight in SurfaceNet (wo occlusion-aware); ★ curve depicts the proposed trainable relative weight in SurfaceNet+. (c): Evaluation of the multi-scale feature aggregation strategy that improves the completeness under different number of view pairs.

lation. Fig. 5(c) further depicts the relationship between sparsity and the average baseline angle over all the models in the ETH3D low-resolution training set. The average baseline angle is far less than $8°$, indicating that the ETH3D dataset may not be suitable for the sparse-MVS benchmark.

## 5 ABLATION STUDY

To investigate the influences of each of the key components in the proposed method, we design an ablation study with respective to the coarse-to-fine fashion (**Multi-scale**) and the trainable occlusion-aware view selection (**View-selection**). For all these studies below, experiments are performed and evaluated on a specific model (model 23) in the DTU dataset because it contains many challenging cases such as complex geometry, textureless regions, and repeating patterns.

In the sparse case, for example, $sparsity = 5$, we quantitatively illustrate the performance gain of the multi-scale fashion in Fig. 10(a), in which we compare few settings: ICCV SurfaceNet [5] (○ curve), ICCV SurfaceNet with the new backbone (● curve) denoted as SurfaceNet in the rest of the paper, SurfaceNet with multi-scale inference (▲ curve), and the proposed trainable occlusion-aware view selection scheme (★ curve). Clearly, from

the comparison of ▲ v.s. ★, we can conclude that the proposed trainable occlusion-aware view selection scheme that is a volume-wise strategy significantly improves both completeness (Recall) and accuracy (Precision).

### 5.1 Multi-scale Mechanism

Fig. 11(a) shows the predictions of the various scale levels. Note that the volume-wise occlusion detection is turned off. Fig. 11(b) contains the result without using the coarse-to-fine mechanism, which is the same as SurfaceNet [5]. The reference model scanned by laser is placed in Fig. 11(c). In each group, (top) the front view of the reconstruction model and (bottom) the intersection of a horizontal plane (red line) are shown. The top view of the red line is useful to observe the surface thickness, noise level, and completeness.

Comparing (a) and (b), it is obvious that the method with the coarse-to-fine mechanism leads to higher precision at the texture area and complex geometry region. Although (b) accurately predicts the results at some complex regions, it suffers from thick surface prediction and floating noise around the repeating pattern regions. The floating noise occurs close to the real surface, because the volume-wise method processes each sub-volume locally and

(a) Multi-scale results of SurfaceNet+        (b) SurfaceNet        (c) Reference Model
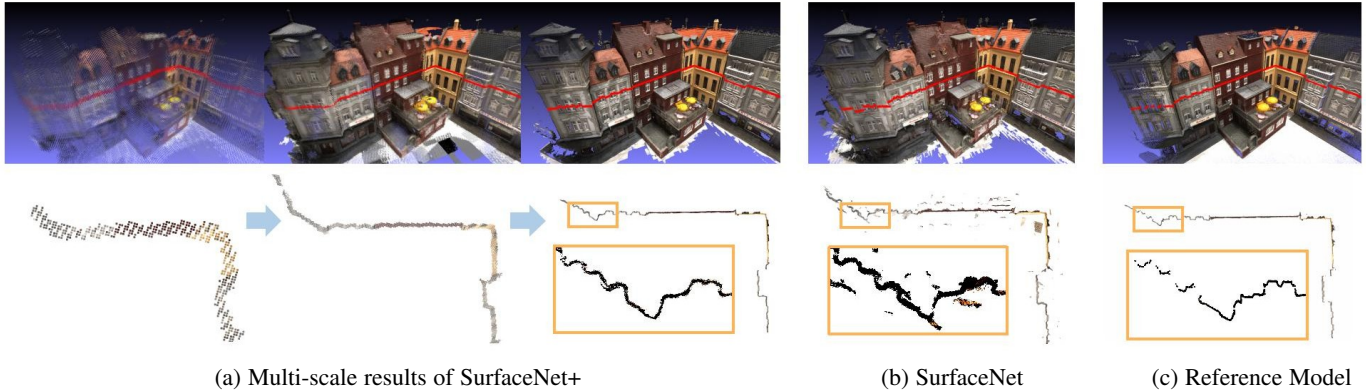
Fig. 11: (a): the predictions of three different scales by the coarse-to-fine mechanism, which gradually refines the geometric information. (b): the reconstructed result without the coarse-to-fine mechanism, i.e., SurfaceNet [5]. (c): the reference model scanned by laser. In each group, we show both (top) the front view of the model and (bottom) the intersection with the red horizontal plane.



(a) Without occlusion detection        (b) **SurfaceNet+(with occlusion detection)**
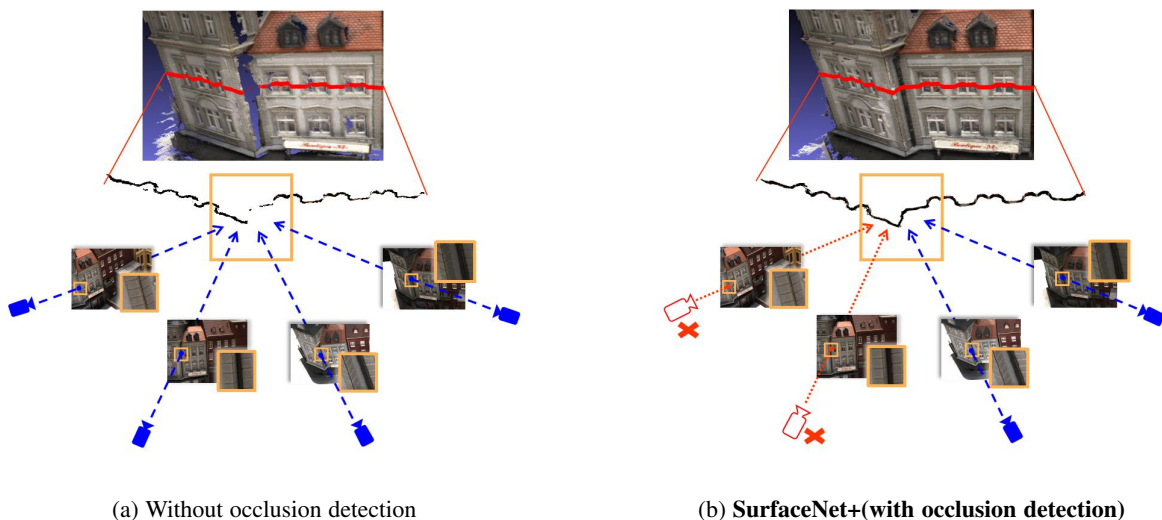
Fig. 12: Qualitative analysis on occlusion detection module. Top: predicted 3D model with selected region (red). Middle: top view of the selected region. Bottom: illustration of the selected (red) / rejected (blue) views. (a): the algorithm without occlusion detection leads to large hole around complex geometry, bounded by a yellow square. (b): occlusion-aware view selection is performed by considering geometric prior and significantly improves the recall (completeness).

individually without global prior to filter out the floating noise. In contrast, the coarse-to-fine mechanism is helpful to gradually reject the empty space and to refine the geometric prediction.

In the sparse case, for example, $n = 5$, the multi-scale mechanism dramatically improves precision if we compare the round-curve and the triangle-curve in Fig. 10(a). Apparently the triangle-curve is a shifted version of the round-curve towards the direction for better precision with constant recall.

**Feature aggregation.** To give the network more global context, we try to use some features coming from the coarser level of the network so that the coarse level is used to not only decide on the visibility/occlusions, but also provide additional feature contents. We study the advantages and disadvantages of this multi-scale inference architectures and report the results in Fig. 10(c). It can be shown that the multi-scale feature aggregation scheme ($\star$) improves the completeness (Recall) of the results by providing the global context. However, when there are few numbers of view pairs, e.g., less than 6 view pairs, the multi-scale aggregation worsens the accuracy (Precision) of the prediction.

The reason is that the volume-wise surface prediction relies on multiple pairs of the unprojected sub-volumes, and in the coarse-to-fine procedure, the selected views may be updated based on the geometric priors under different scales. So that when the multi-scale scheme aggregates the features from different view pairs, the global context may become less useful and leads to worse accuracy (Precision).

## 5.2 Occlusion Detection

To analyze the qualitative impact of occlusion-aware view selection, the comparison experiment is set based on the result using the multi-scale mechanism. For better visualization, we only probe the occlusion-aware view selection in the final multi-scale stage. The results with and without occlusion-aware view selection are shown in Fig. 12, which contains the front views and the intersection (the red line shown in the model) of the results accompanied by different camera views.

Note that SurfaceNet+ (with View-selection) has higher recall output, especially around the corner of the reconstructed house

model (shown in the orange box of the intersection). The gap lies in different views selected by each method. Both methods use patch image (bottom right corner of the picture) to select valid views (the four views shown in the bottom of the figure). Yet the left two views are blocked by the surface, which means only the right two views can provide useful patch information for reconstruction. The occluded views reduced the output weight under the correct views; therefore, incomplete prediction occurred in complex geometry regions without occlusion-aware view selection. In SurfaceNet+, the rejected occluded views (shown in red) are detected by the projection rays combined with the output point cloud in the previous stage mentioned in subsection3.2. It is worth noting that these occluded views are extremely hard to detect using only image patches. These patches are similar to each other, so it is difficult to infer the relative position relationship among them in the absence of three-dimension prior.

In Fig. 10(b), to further demonstrate the benefit from an explicit "relative weight (with occlusion-aware)" (★ curve), we investigate the setting "relative weight (without occlusion-aware)" (▲ curve) and the setting "without relative weight" (□ curve). Enabling the "relative weight (with occlusion-aware)" significantly improves Recall (Completeness) of the reconstructed model, indicating the effectiveness of the proposed trainable occlusion-aware view selection scheme.

Additionally, in Fig. 10(a), we evaluate the proposed end-to-end trainable occlusion-aware view selection scheme (trainable, ★ curve) versus the heuristic view selection method (heuristic, ■ curve). Note that both of them share the same backbone network structure and the multiscale fusion strategy, while the only difference is the view selection module. As we can see, the proposed end-to-end trainable occlusion-aware view selection scheme significantly boosts the completeness (Recall) of the reconstruction model.

**SurfaceNet [5].** For fair comparison, in Fig. 10(a) we also show the performance difference between the ICCV SurfaceNet [5] and the modified SurfaceNet with the new backbone, where the only modification on the ICCV SurfaceNet is the network structure that SurfaceNet+ is using (Fig. 4). It is worth noting the relative position changes of each curve. There is a clear shift downward after adding the proposed trainable occlusion-aware view selection(**View-selection**). This indicates the better recall with comparable precision. Overall, the gain achieved by SurfaceNet+ over SurfaceNet has NO relationship with the backbone network adopted; instead, it is benefited from the proposed multiscale pipeline and novel view selection strategy.

### 5.3 Discussion

**Hyperparameters.** The number of view pairs $N_v$ is also critical for the algorithm. Too few view pairs may lead to noisy and inaccurate reconstruction, while too many in sparse-MVS lead to incomplete (low recall) prediction. The trade-off of $N_v = 3$ achieves the best overall performance as indicated in the figure.

To further analyze the effect of occlusion-aware view selection, we experiment with different occlusion parameters $\alpha$ at different sparsities with a fixed view pair number $N_v = 3$. The recall gain is counted by the recall improvement based on the method without occlusion detection.

As shown in Fig. 13, the gain increases as the sparsity grows. The reason lies in that when sparsity increases, a growing baseline angle and fewer view pairs lead to a lower percentage of
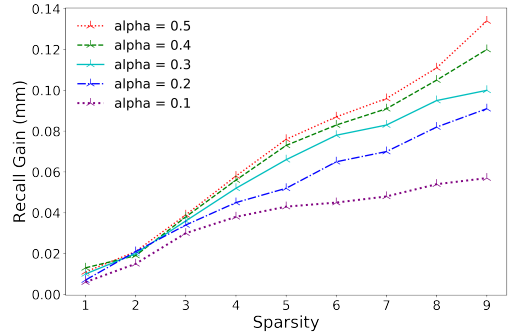


Fig. 13: Recall gain w.r.t. $\alpha$(occlusion parameter) at different sparsity. The gain is counted by the recall improvement based on the method without occlusion detection.

non-occluded views. Therefore, lower weight on occluded views controlled by alpha has increased benefit on larger sparsity.

| | NO. of sub-volumes | | Speed up | | |
|---|---|---|---|---|---|
| | SurfaceNet | SurfaceNet+ | mean | max | min |
| DTU | 140,608 | 12,320 | **11X** | 23X | 7X |
| T&T | 158,992 | 15,892 | **15X** | 33X | 11X |

TABLE 3: Efficiency comparison of proposed method with and without volume selection. Where we set the resolution as 4(mm) in the Tanks and Temple dataset (T&T) and 0.03(mm) in the DTU dataset and compare the average cubic number and its mean and maximum speed up ratio for each model.
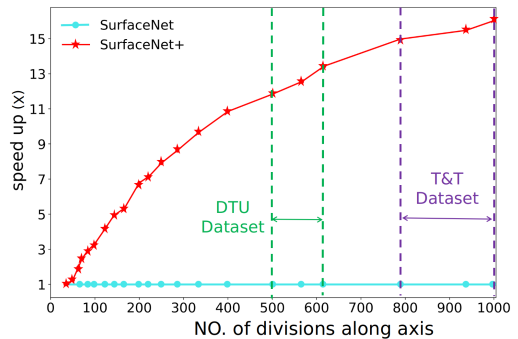


Fig. 14: Speed up ratio with the change of resolution. Note how coarse to fine mechanism leads to efficient representation compared to SurfaceNet [5]. With the finer of reconstruction, speed up ratio grows dramatically.

**Efficiency.** To evaluate the efficiency brought by the coarse-to-fine mechanism, we measure the speed of the algorithm using the total sampled volumes. Specifically, we count each number of cubes sampled by the algorithm for both methods in all the models on the DTU [10] evaluation set and Tanks and Temples [23] 'Intermediate' set. We set the whole reconstruction scene as a cubic box with length $l_{scene} = 400(mm)$ in the DTU dataset and the final voxel resolution $r = 0.3(mm)$. Each volume forms a tensor of size $s \times s \times s$ and we set $s = 64$. We use all the cubes for reconstruction in SurfaceNet [5], a three-stage coarse-to-fine pipeline for SurfaceNet+. The settings in Tanks and Temples are equal to DTU except that we set $l_{scene} = 400(mm)$ and $r = 2(mm)$. The left part of Table 3 shows the average sub-volumes used for reconstruction, and the right part shows the speed up multiple brought by the coarse-to-fine mechanism. We value

the average multiple as the ratio between two methods. The mean, maximum and minimum multiple show that the volume selection mechanism can achieve more than 10 times higher efficiency on both datasets.

To better understand the efficiency promotion, in Fig. 14, we visualize the speed up ratio as the scale of the relative resolution $r_{relative} = \frac{l_{scene}}{r}$ . Note how the coarse-to-fine mechanism leads to efficient representation compared to SurfaceNet. At low relative resolution, the ratio is near 1 due to the nearly dense sampling based on the coarse prediction. Yet with the finer reconstruction, the speedup ratio grows dramatically because the finer prediction leads to a higher percentage of empty sub-volumes.
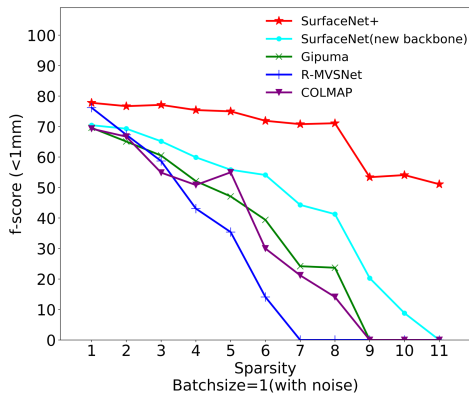


Fig. 15: Evaluation of the sparse MVS benchmark using the **noisy camera poses** (estimated by SfM [56]),

**Noisy camera poses.** The camera poses used in our previous experiments are given by the public datasets, which are estimated by the registration of laser scans (denoted as GT camera pose). While in practice, the camera poses may be computed through the sparse set of views, which inevitably suffers noise (denoted as noisy camera pose). To evaluate how the noisy camera pose affects the performance of SurfaceNet+, we adopt the structure-of-motion SfM [56] along with the sparse set of views to obtain the noisy camera pose. As expected, using Noisy camera poses (Fig. 15) degrades the performance of MVS methods that using GT camera poses (Fig. 6), where the f-score drops.

We examine the f-score degradation between Fig. 15 and Fig. 6, where the image-wise view selection scheme, used in Gipuma [6] and R-MVSNet [7], is more sensitive to the camera pose noise, especially under massive sparsity levels. In contrast, the pixel-wise (COLMAP [9]) and volume-wise (SurfaceNet [5] and SurfaceNet+) view selection strategy is relatively more robust to camera pose noise. The reason is that the camera pose noise will introduce an inhomogeneous shift of the photo-consistent matches, so that the pixel-wise and volume-wise view selection can adaptively choose the relatively better views based on the photometric consistency despite the noisy camera pose. In contrast, the image-wise view selection leads to matching the correspondence only on the pre-selected views, which no longer be the best views for a large proportion of pixels or sub-volumes if the camera pose noise is considered.

## 6  CONCLUSION

As sparser sensation is more practical and more cost-efficient, instead of only focusing on dense MVS setup, we propose a comprehensive analysis on sparse-MVS under various observation sparsities. The proposed leader-board calls for more attention and effort from the community to the sparse-MVS problem, since the state-of-the-art depth-fusion methods significantly perform worse as the baseline angle get larger in the sparser setting. As another line of the solution, we propose a volumetric method, SurfaceNet+, to handle sparse-MVS by introducing the novel occlusion-aware view selection scheme as well as the multi-scale strategy. Consequently, the experiments demonstrate the tremendous performance gap between SurfaceNet+ and the recent methods in terms of precision and recall. Under the extreme sparse-MVS settings in two datasets, where existing methods can only return very few points, SurfaceNet+ still works as well as in the dense MVS setting.

**Limitations.** (1) Ideally, for a simple geometric region, each piece of a surface in sub-volume should be effectively reconstructed only using ONE view pair with large baseline angle, i.e., $N_v = 1$. However, due to various of shading and lighting conditions, the colorization of the 3D model gets more challenging by using less number of views. (2) Furthermore, even though the scanned models in the MVS datasets are large-scale scene, it will be challenging for SurfaceNet+ to effectively and efficiently reconstruct a city-level 3D model. (3) Last but not least, despite of the great generalization ability of the learnt model, it still requires dozens of laser-scanned 3D model for supervision. That significantly limits the application scenarios, such as astro-observation and multi-view microscopic observation, where rare supervision signal can be captured.

## REFERENCES

[1] K. Yucer, C. Kim, A. Sorkine-Hornung, and O. Sorkine-Hornung, "Depth from gradients in dense light fields for object reconstruction," in *2016 Fourth International Conference on 3D Vision (3DV)*. IEEE, 2016, pp. 249–257.

[2] G. Wu, Y. Liu, L. Fang, Q. Dai, and T. Chai, "Light field reconstruction using convolutional network on epi and extended applications," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 7, pp. 1681–1694, 2018.

[3] T. E. Bishop, S. Zanetti, and P. Favaro, "Light field superresolution," in *2009 IEEE International Conference on Computational Photography (ICCP)*. IEEE, 2009, pp. 1–9.

[4] J. Chung, G. W. Martinez, K. C. Lencioni, S. R. Sadda, and C. Yang, "Computational aberration compensation by coded-aperture-based correction of aberration obtained from optical fourier coding and blur estimation," *Optica*, vol. 6, no. 5, pp. 647–661, 2019.

[5] M. Ji, J. Gall, H. Zheng, Y. Liu, and L. Fang, "Surfacenet: An end-to-end 3d neural network for multiview stereopsis," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2307–2315.

[6] S. Galliani, K. Lasinger, and K. Schindler, "Massively parallel multiview stereopsis by surface normal diffusion," in *IEEE International Conference on Computer Vision*, 2015, pp. 873–881.

[7] Y. Yao, Z. Luo, S. Li, T. Shen, T. Fang, and L. Quan, "Recurrent mvsnet for high-resolution multi-view stereo depth inference," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5525–5534.

[8] R. Chen, S. Han, J. Xu, and H. Su, "Point-based multi-view stereo network," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.

[9] J. L. Schönberger, E. Zheng, M. Pollefeys, and J.-M. Frahm, "Pixelwise view selection for unstructured multi-view stereo," in *European Conference on Computer Vision (ECCV)*, 2016.

[10] H. Aanæs, R. R. Jensen, G. Vogiatzis, E. Tola, and A. B. Dahl, "Large-scale data for multiple-view stereopsis," *International Journal of Computer Vision*, pp. 1–16, 2016.

[11] W. Chen, Z. Fu, D. Yang, and J. Deng, "Single-image depth perception in the wild," in *Advances in neural information processing systems*, 2016, pp. 730–738.

[12] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 270–279.

[13] F. Liu, C. Shen, and G. Lin, "Deep convolutional neural fields for depth estimation from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 5162–5170.

[14] Z. Huang, T. Li, W. Chen, Y. Zhao, J. Xing, C. LeGendre, L. Luo, C. Ma, and H. Li, "Deep volumetric video from very sparse multi-view performance capture," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 336–354.

[15] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 4460–4470.

[16] A. Kar, C. Häne, and J. Malik, "Learning a multi-view stereo machine," in *Advances in neural information processing systems*, 2017, pp. 365–376.

[17] 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.

[18] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, "Mvsnet: Depth inference for unstructured multi-view stereo," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 767–783.

[19] X. Gu, Z. Fan, S. Zhu, Z. Dai, F. Tan, and P. Tan, "Cascade cost volume for high-resolution multi-view stereo and stereo matching," *arXiv preprint arXiv:1912.06378*, 2019.

[20] C. Barnes, E. Shechtman, A. Finkelstein, and D. B. Goldman, "Patchmatch: A randomized correspondence algorithm for structural image editing," *ACM Trans. Graph*, vol. 28, no. 3, pp. 24–1, 2009.

[21] J. Pang, O. C. Au, Y. Yamashita, Y. Ling, Y. Guo, and J. Zeng, "Self-similarity-based image colorization," in *IEEE International Conference on Image Processing*, 2014, pp. 4687–4691.

[22] A. Zheng, Y. Yuan, S. P. Jaiswal, and O. C. Au, "Motion estimation via hierarchical block matching and graph cut," in *IEEE International Conference on Image Processing*, 2015, pp. 4371–4375.

[23] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and temples: Benchmarking large-scale scene reconstruction," *ACM Transactions on Graphics*, vol. 36, no. 4, 2017.

[24] M. Lhuillier and L. Quan, "A quasi-dense approach to surface reconstruction from uncalibrated images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 3, pp. 418–433, March 2005.

[25] Y. Furukawa and J. Ponce, "Accurate, dense, and robust multiview stereopsis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 8, pp. 1362–1376, Aug 2010.

[26] E. Tola, C. Strecha, and P. Fua, "Efficient large-scale multi-view stereo for ultra high-resolution image sets," *Machine Vision and Applications*, pp. 1–18, 2012.

[27] N. D. Campbell, G. Vogiatzis, C. Hernández, and R. Cipolla, "Using multiple hypotheses to improve depth-maps for multi-view stereo," in *European Conference on Computer Vision*, 2008, pp. 766–779.

[28] P. Merrell, A. Akbarzadeh, L. Wang, P. Mordohai, J. Frahm, R. Yang, D. Nistér, and M. Pollefeys, "Real-time visibility-based fusion of depth maps," in *IEEE International Conference on Computer Vision*, 2007, pp. 1–8.

[29] C. Zach, "Fast and high quality fusion of depth maps," 01 2008.

[30] V. Lempitsky and Y. Boykov, "Global optimization for shape fitting," in *2007 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2007, pp. 1–8.

[31] B. Curless and M. Levoy, "A volumetric method for building complex models from range images," in *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996, pp. 303–312.

[32] G. Riegler, A. O. Ulusoy, H. Bischof, and A. Geiger, "Octnetfusion: Learning depth fusion from data," in *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, pp. 57–66.

[33] K. N. Kutulakos and S. M. Seitz, "A theory of shape by space carving," in *IEEE International Conference on Computer Vision*, vol. 1, 1999, pp. 307–314.

[34] S. M. Seitz and C. R. Dyer, "Photorealistic scene reconstruction by voxel coloring," *International Journal of Computer Vision*, vol. 35, no. 2, pp. 151–173, 1999.

[35] A. Kar, C. Häne, and J. Malik, "Learning a multi-view stereo machine," 2017.

[36] M. Jancosek and T. Pajdla, "Multi-view reconstruction preserving weakly-supported surfaces," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2011, pp. 3121–3128.

[37] 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.

[38] M. Kazhdan, M. Bolitho, and H. Hoppe, "Poisson surface reconstruction," in *Proceedings of the fourth Eurographics symposium on Geometry processing*, vol. 7, 2006.

[39] C. Häne, S. Tulsiani, and J. Malik, "Hierarchical surface prediction for 3d object reconstruction," in *2017 International Conference on 3D Vision (3DV)*. IEEE, 2017, pp. 412–420.

[40] X. Han, T. Leung, Y. Jia, R. Sukthankar, and A. C. Berg, "Matchnet: Unifying feature and metric learning for patch-based matching," in *CVPR*, 2015.

[41] A. Seki and M. Pollefeys, "Sgm-nets: Semi-global matching with neural networks," in *CVPR*, 2017.

[42] P. Knöbelreiter, C. Reinbacher, A. Shekhovtsov, and T. Pock, "End-to-End Training of Hybrid CNN-CRF Models for Stereo," in *2017 Computer Vision and Pattern Recognition (CVPR)*, 2017.

[43] S. Galliani and K. Schindler, "Just look at the image: Viewpoint-specific surface normal prediction for improved multi-view reconstruction," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5479–5487.

[44] A. Kendall, H. Martirosyan, S. Dasgupta, P. Henry, R. Kennedy, A. Bachrach, and A. Bry, "End-to-end learning of geometry and context for deep stereo regression," *CoRR*, vol. abs/1703.04309, 2017.

[45] L. Xu, L. Hou, O. C. Au, W. Sun, X. Zhang, and Y. Guo, "A novel ray-space based view generation algorithm via radon transform," *3D Research*, vol. 4, no. 2, p. 1, 2013.

[46] Y. Furukawa and J. Ponce, "Accurate, dense, and robust multiview stereopsis," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 8, pp. 1362–1376, 2010.

[47] M. Jancosek and T. Pajdla, "Multi-view reconstruction preserving weakly-supported surfaces," 07 2011, pp. 3121 – 3128.

[48] M. Goesele, N. Snavely, B. Curless, H. Hoppe, and S. Seitz, "Multi-view stereo for community photo collections," 11 2007, pp. 1–8.

[49] T. Schops, J. L. Schonberger, S. Galliani, T. Sattler, K. Schindler, M. Pollefeys, and A. Geiger, "A multi-view stereo benchmark with high-resolution images and multi-camera videos," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3260–3269.

[50] R. Jensen, A. Dahl, G. Vogiatzis, E. Tola, and H. Aanæs, "Large scale multi-view stereopsis evaluation," in *2014 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2014, pp. 406–413.

[51] Q. Xu and W. Tao, "Multi-scale geometric consistency guided multi-view stereo," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5483–5492.

[52] S. Cheng, Z. Xu, S. Zhu, Z. Li, L. E. Li, R. Ramamoorthi, and H. Su, "Deep stereo using adaptive thin volume representation with uncertainty awareness," *arXiv preprint arXiv:1911.12012*, 2019.

[53] J. Liao, Y. Fu, Q. Yan, and C. Xiao, "Pyramid multi-view stereo with local consistency," in *Computer Graphics Forum*, vol. 38, no. 7. Wiley Online Library, 2019, pp. 335–346.

[54] R. Chen, S. Han, J. Xu *et al.*, "Visibility-aware point-based multi-view stereo network," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[55] R. Chen, S. Han, J. Xu, and H. Su, "Point-based multi-view stereo network," in *The IEEE International Conference on Computer Vision (ICCV)*, 2019.

[56] J. L. Schönberger and J.-M. Frahm, "Structure-from-motion revisited," in *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

**Mengqi Ji** is currently a postdoc in Tsinghua University. He received Ph.D / M.Sc from the Hong Kong University of Science and Technology in 2019 / 2013, and B.E. from University of Science and Technology Beijing in 2012. His research interests include 3D vision and computational photography.

**Jinzhi Zhang** is currently a master student in Tsinghua-Berkeley Shenzhen Institute (TBSI), Tsinghua University. He received B.E. from Huazhong University of Science and Technology in 2019. His research interest is 3D vision.

**Qionghai Dai** is a professor in the Department of Automation, and an adjunct professor in the School of Life Science, Tsinghua University. Dr. Dai is the academician of Chinese Academy of Engineering. His research interests include computational photography, brain science, and artificial intelligence. Dr. Dai is currently IEEE Senior Member, serving as Associate Editor for IEEE TIP.

**Lu Fang** is currently an Associate Professor in Tsinghua University. She received Ph.D from the Hong Kong Univ. of Science and Technology in 2011, and B.E. from Univ. of Science and Technology of China in 2007. Her research interests include computational photography and 3D vision. Dr. Fang used to receive Best Student Paper Award in ICME 2017, Finalist of Best Paper Award in ICME 2017 and ICME 2011. Dr. Fang is currently IEEE Senior Member.