# M$^3$VSNet: Unsupervised Multi-metric Multi-view Stereo Network

Baichuan Huang, Hongwei Yi, Can Huang, Yijia He, Jingbin Liu, Xiao Liu
Wuhan University, Peking University, Megvii Technology Limited

## ABSTRACT

The present Multi-view stereo (MVS) methods with supervised learning-based networks have an impressive performance comparing with traditional MVS methods. However, the ground-truth depth maps for training are hard to be obtained and are within limited kinds of scenarios. In this paper, we propose a novel unsupervised multi-metric MVS network, named M$^3$VSNet, for dense point cloud reconstruction without any supervision. To improve the robustness and completeness of point cloud reconstruction, we propose a novel multi-metric loss function that combines pixel-wise and feature-wise loss function to learn the inherent constraints from different perspectives of matching correspondences. Besides, we also incorporate the normal-depth consistency in the 3D point cloud format to improve the accuracy and continuity of the estimated depth maps. Experimental results show that M$^3$VSNet establishes the state-of-the-arts unsupervised method and achieves comparable performance with previous supervised MVSNet on the *DTU* dataset and demonstrates the powerful generalization ability on the *Tanks & Temples* benchmark with effective improvement.

## CCS CONCEPTS

• **Computing methodologies** → **Reconstruction**.

## 1 INTRODUCTION

Multi-view stereo (MVS) aims to reconstruct the 3D dense point cloud from multi-view images [6, 9], which has various applications in augmented reality, virtual reality and robotics, etc. [24]. Big progress has been made in the dense reconstruction with traditional methods through the hand-crafted features (e.g. NCC) to calculate the matching correspondences [10–12, 23, 28]. Though, the efficient and robust methods of MVS in the large-scale environments are still the challenging tasks [24]. Recently, deep learning is introduced to relieve this limitation. The supervised learning-based MVS methods achieve significant progress especially improving the efficiency and completeness of dense point cloud reconstruction [15, 27, 31,

32]. These learning-based methods learn and infer the information to handle matching ambiguity which is hard to be obtained by stereo correspondences. However, these supervised learning-based methods strongly depend on the training datasets with ground-truth depth maps, which have limited kinds of scenarios and are not easy to be available. Thus it is a big hurdle and may lead to bad generalization ability in different complex scenarios [7, 13, 30, 31]. Furthermore, the robustness and completeness of dense point cloud reconstruction still have a lot of room to be improved. Current unsupervised learning-based methods are mainly based on the pixel-wise level, which will cause incorrect matching correspondences with low robustness [30][31]. Because for two identical images, the difference will be huge as long as pixel offset from the perspective of pixel level. However, they are almost the same from the perspective of perception such as feature level. In addition, the human visual system perceives the surrounding world depending on the object features rather than a single image pixel [4]. Therefore, in order to improve the robustness and completeness of unsupervised learning-based MVS, it drives us to consider the similarity on object features.

In this paper [1], we propose a novel unsupervised multi-metric MVS network, named M$^3$VSNet as shown in figure 1, which could infer the depth maps for dense point cloud reconstruction even in non-ideal environments. Most importantly, we propose a novel multi-metric loss function, namely pixel-wise and feature-wise loss function. The key insight is that the human visual system perceives the surrounding world by the object features [4]. In terms of this loss function, both the photometric and geometric matching consistency can be well guaranteed, which is more accurate and robust compared with the only photometric constraints used in MVSNet [31]. Specifically, we introduce the multi-scale feature maps from the pre-trained VGG16 network as vital clues in the feature-wise loss. Low-level feature representations learn more texture details while high-level features learn semantic information with a large receptive field. Different level features are the representations of different receptive fields. By aggregating multi-scale features, our proposed M$^3$VSNet can consider both the low-level image texture and the high-level semantic information. Therefore, the network can well improve the robustness and accuracy of matching correspondences. Compared with the network only using pixel-wise loss which performs mismatch errors in some challenging scenarios such as textureless, mirror effect or reflection and texture repeat areas [24, 30, 31], M$^3$VSNet can improve the robustness by considering the similarity between the multi-scale semantic features.

Besides, in order to further improve the performance of the estimated depth maps, we incorporate the normal-depth consistency in the world coordinate space to constraint the local surface tangent obtained from the estimated depth maps to be orthogonal to

---

[1]Can Huang is the corresponding author

the calculated normal. This regularization will improve the accuracy and continuity of the estimated depth maps. Moreover, we utilize the multi-scale pyramid feature aggregation to construct the 3D cost volume with more contextual information to improve the robustness and accuracy of feature correspondences.

Our main contributions are summarized as below:

- We propose a novel multi-metric unsupervised network, which can work even in non-ideal environments, for multi-view stereo without any ground-truth 3D training data.
- we propose a novel multi-metric loss function that considers different perspectives of matching correspondences beyond pixel value. Besides, we incorporate the normal-depth consistency in the 3D point cloud format to improve the accuracy and continuity of the estimated depth maps.
- Extensive experiments demonstrate that our proposed $M^3$VSNet outperforms the previous state-of-the-art unsupervised methods and achieves comparable performance with original MVSNet on the *DTU* dataset and shows the excellent generalization ability on the *Tanks & Temples* benchmark with effective improvement.

## 2 RELATED WORK

### 2.1 Traditional MVS

Many traditional methods have been proposed in this field such as voxel-based method [25], feature points diffusion [10] and the fusion of estimated depth maps [3]. First of all, the voxel-based method consumes many computing resources and its accuracy depends on the resolution of voxel mainly [15]. Secondly, the blank area may seriously suffer from the textureless problem in the method of feature points diffusion. Thirdly, the most used method is the fusion of inferred depth maps, which gets the depth maps and then fuses all the depth maps together to the final point cloud [5]. Besides, many methods of improvement have been proposed. Silvano [11] formulates the patch matches in 3D space and the progress can be massively parallelized and delivered. Johannes [23] estimates the depth and normal maps synchronously and uses photometric and geometric priors to refine the image-based depth and normal fusion. Though, the robustness needs to be improved when dealing with non-ideal environments such as textureless or texture repeat areas and no-Lambert surfaces.

### 2.2 Depth Estimation

The fusion of estimated depth maps can decouple the reconstruction into depth estimation and depth fusion. Depth estimation with monocular video and binocular image pairs has many similarities with the multi-view stereo here [19]. But there are exactly some differences between them. Monocular video [34] lacks the real scale of the depth actually and binocular image pairs always need to rectify the parallel two images [8]. In this case, only the disparity needs to be inferred without considering the intrinsic and extrinsic of the camera. As for multi-view stereo, the input is the arbitrary number of pictures. What's more, the transformation among these positions should be taken into consideration as a whole [31]. Other obstacles such as multi-view occlusion and consistency [7] raise the bar for depth estimation of multi-view stereo than that of monocular video and binocular image pairs.

### 2.3 Supervised Learning MVS

Since Yao Yao proposed MVSNet in 2018 [31], many supervised networks based on MVSNet have been proposed. To reduce GPU memory consumption, Yao Yao introduces R-MVSNet with the help of gated recurrent unit [32]. Gu uses the concept of the cascade to shrink the cost volume [13]. Yi introduced two new self-adaptive view aggregation with pyramid multi-scale images to enhance the point cloud in textureless regions [33]. Luo utilizes the plane-sweep volumes with isotropic and anisotropic 3D convolutions to get better results [21]. In this kind of task, cost volume and 3D regularization are highly memory-consuming. More importantly, the ground-truth depth maps are derived from heavy labor.

### 2.4 Unsupervised Learning MVS

The unsupervised network utilizes the photometric and geometric constraints to learn the depth by itself, which relief the complicated artificial markers for ground-truth depth maps. Many works explore unsupervised learning in monocular video and binocular image pairs. Reza [22] presents the unsupervised learning method for depth and ego-motion from monocular video. The paper uses image reconstruction loss, 3D point cloud alignment loss and additional image-based loss. Being similar to unsupervised learning in monocular video and binocular image pairs [2], the losses of MVS are also based on photometric and geometric consistency. Dai [7] predicts the depth maps for all views simultaneously in a symmetric way. In the stage, cross-view photometric and geometric consistency can be guaranteed. But this method consumes a lot of GPU memory. Additionally, Tejas [17] proposes the simplified network and traditional loss designation but an unsatisfied result. Efforts are worthy to be paid in this direction.

## 3 $M^3$VSNET

In this section, we introduce our proposed $M^3$VSNet in detail. We firstly describe the network architecture in section 3.1 to generate initial depth map, then illustrate the normal-depth consistency in section 3.2 to refine it in consideration of the orthogonality between normal and local surface tangent. Finally, our proposed novel multi-metric loss in section 3.3 is introduced by considering different perspectives of matching correspondences to improve the robustness and completeness of point cloud reconstruction.

### 3.1 Network Architecture

The basic architecture of our proposed $M^3$VSNet consists of three parts, namely pyramid feature aggregation, variance-based cost volume generation and 3D U-Net regularization, as shown in figure 1. The pyramid feature aggregation extracts features from low-level to high-level representations with more contextual information. Then the same variance-based cost volume generation and 3D U-Net regularization as MVSNet [31] are used to generate the initial depth map. The advance architecture of $M^3$VSNet consists of two parts, namely normal-depth consistency and multi-metric loss. After generating the initial depth map, we incorporate the novel normal-depth consistency to refine it in consideration of the orthogonality between normal and local surface tangent. More importantly, we construct multi-metric loss, which consists of pixel-wise loss and feature-wise loss. We will briefly describe each module in the following parts.
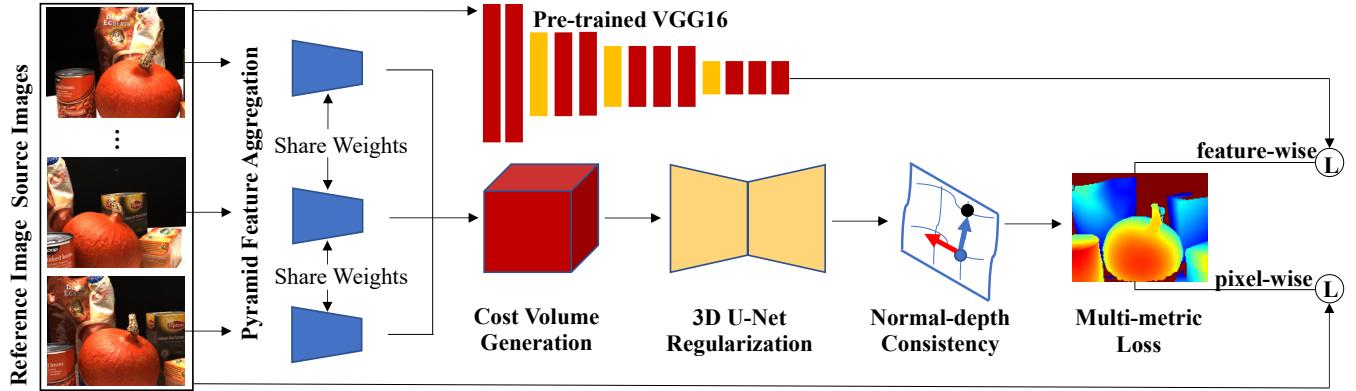
**Figure 1: The architecture of our proposed M³VSNet. It contains five components: pyramid feature aggregation, variance-based cost volume generation, 3D U-Net regularization, normal-depth consistency and multi-metric loss function.**
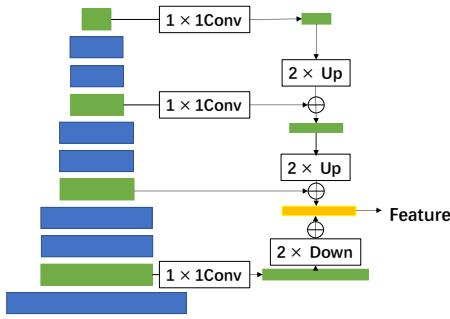


**Figure 2: The illustration of pyramid feature aggregation.**



**Figure 3: The illustration of normal from depth**

*3.1.1 Pyramid Feature Aggregation.* In previous supervised learning-based network such as MVSNet [31], only the 1/4 feature is adopted (1/4 represents a quarter of the size of the original reference image). 1/4 feature lacks multi-scale context information for matching correspondences. Therefore, we propose to use the pyramid feature aggregation that aggregates different scale features with contextual information of different receptive fields [20]. Figure 2 shows the details of this module. For the input images, the feature extraction network is constructed to extract the aggregated 1/4 feature. In the process of bottom-up, the stride of the layer 3, 6 and 9 is set to 2 to get the four scale features in eleven-layer 2D CNN. Each convolutional layer is followed by the structure of BatchNorm and ReLU. In the process of up-bottom, each level of features is derived from the concatenate by the upsampling of the higher layer and the feature in the same layer with fewer channels. Especially, the 1/2 feature needs to be downsampled to be aggregated into the final 1/4 feature. To reduce the dimension of the final 1/4 feature, the $1 \times 1$ convolution for each concatenation is adopted. At last, we get the final feature with 32 channels, which is an aggregation of contextual information from low-level to high-level representations.

*3.1.2 Cost volume and 3D U-Net regularization.* The construction of variance-based cost volume is based on the differentiable homography warping with the number of different depth hypotheses $D$ in
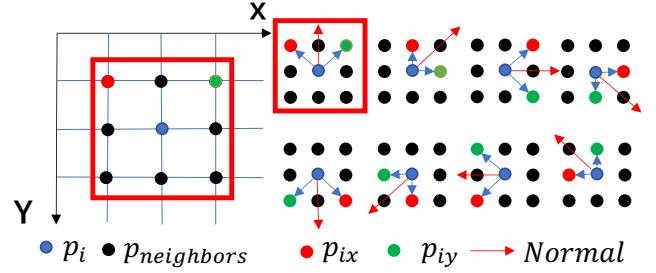
MVSNet [31]. Then 3D U-Net regularization is used to regularize the 3D cost volume, which is simple but effective for aggregating features. At last, the initial depth is derived from the *soft argmin* operation with the probability volume after the regularization.

## 3.2 Normal-depth Consistency

The initial depth still contains some incorrect matching correspondences. Therefore, to improve the quality of estimated depth maps further, we incorporate the normal-depth consistency based on the orthogonality between normal and local surface tangent [30]. The consistency will make the depth more reasonable in 3D space. Normal-depth consistency can be divided into two steps. Firstly, the normal should be calculated by the depth with the orthogonality. Then the refined depth can be inferred by the normal and initial depth according to the projection relationship. This module cooperating with 3D U-Net will refine the depth in 2D and 3D space jointly, which improves the accuracy and continuity of depth.

As shown in figure 3, eight neighbors are selected to infer the normal of the central point. Due to the orthogonality, the operation of cross-product is used. For each central point $p_i$, one set of the neighbors can be recognized as $p_{ix}$ and $p_{iy}$. If the depth $Z_i$ of $p_i$ and the intrinsics $K$ of camera are known, the normal $\widetilde{N_i}$ can be calculated as below:
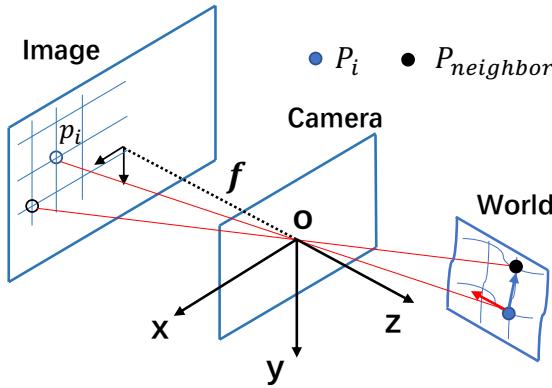
$$P_i = K^{-1} Z_i p_i \tag{1}$$

**Figure 4: The illustration of depth from normal**

$$\widetilde{N_i} = \overrightarrow{P_iP_{ix}} \times \overrightarrow{P_iP_{iy}} \tag{2}$$

To add the credibility of final normal estimation $N_i$, mean cross-product for eight neighbors can be presented as below:

$$N_i = \frac{1}{8}\sum_{i=1}^{8}(\widetilde{N_i}) \tag{3}$$

The final refined depth maps can be available when the normal and initial depth maps are provided. In figure 4, for each pixel $p_i(x_i, y_i)$, the depth of the neighbor $p_{neighbor}$ should be refined. Their corresponding 3D points are $P_i$ and $P_{neighbor}$. The normal of $P_i$ is $\overrightarrow{N_i}(n_x, n_y, n_z)$. The depth of $P_i$ is $Z_i$ and the depth of $P_{neighbor}$ is $Z_{neighbor}$. We can get the equation $\overrightarrow{N} \perp \overrightarrow{P_iP_{neighbor}}$. The relationship is apparently reasonable due to the orthogonality and surface consistency in the local surface. In summary, the depth $Z_{neighbor}$ of the neighbors can be inferred by the depth and normal of the central point.

$$(K^{-1}Z_ip_i - K^{-1}Z_{neighbor}p_{neighbor})\begin{bmatrix} n_x \\ n_y \\ n_z \end{bmatrix} = 0 \tag{4}$$

For the refined depth, eight neighbors are also taken into consideration. Considering the discontinuity of normal in some edge or irregular surface, the weight $w_i$ for the reference image $I_i$ is introduced to make depth more conforming to geometric consistency. The weight is defined as below:

$$w_i = e^{-\alpha_1|\nabla I_i|} \tag{5}$$

The weight $w_i$ depends on the gradient between $p_i$ and $p_{neighbor}$, which means that the bigger gradient represents the less reliability of the refined depth. In view of the eight neighbors, the final refined depth $\widetilde{Z}_{neighbor}$ is a combination of the weighted sum of eight different directions. The final refined depth is the result of regularization in 3D space, which improves the accuracy and continuity of the estimated depth maps.

$$\widetilde{Z}_{neighbor} = \sum_{i=1}^{8} w_i' Z_{neighbor} \tag{6}$$

$$w_i' == \frac{w_i}{\sum_{i=1}^{8} w_i} \tag{7}$$

## 3.3 Multi-metric Loss

We propose a novel multi-metric loss function by considering different perspectives of matching in feature correspondence beyond pixel, which is quite crucial and effective. The pixel-wise loss can guarantee the matching correspondences with more texture details and the feature-wise loss can make use of the semantic information.

The key idea embodied in multi-metric loss function is the photometric consistency crossing multi-views [3]. Given the reference image $I_{ref}$ and source image $I_{src}$, the corresponding intrinsic parameters are represented as $K_{ref}$ and $K_{src}$. Besides, the extrinsic from $I_{ref}$ to $I_{src}$ is represented as $T$. For the pixel $p_i(x_i, y_i)$ in $I_{ref}$, the corresponding pixel $p_i'(x_i', y_i')$ in $I_{src}$ can be calculated as:

$$p_i' = KT(K^{-1}\widetilde{Z}_ip_i) \tag{8}$$

The overlapping area, named $I_{src}'$, from reference image $I_{ref}$ to source image $I_{src}$ can be sampled using the bilinear interpolation.

$$I_{src}' = I_{src}(p_i') \tag{9}$$

For the occlusion area, the value of the pixel in $I_{src}'$ is set to zero. Obviously, the mask $M$ can be obtained when the $p_i$ is projected to the external area of $I_{src}$. Based on the prior constraint, the multi-metric loss function $L$ is formulated as the sum of pixel-wise loss $L_{pixel}$ and feature-wise loss $L_{feature}$.

$$L = \sum(\gamma_1 L_{pixel} + \gamma_2 L_{feature}) \tag{10}$$

*3.3.1 Pixel-wise Loss.* For the pixel-wise loss, we only consider the photometric consistency between the reference image $I_{ref}$ and other source images. There are mainly three parts of this loss function. Firstly, the photometric loss compares the difference of pixel value between $I_{ref}$ and $I_{src}'$. To relieve the influence of lighting changes, the gradient of every pixel is integrated into $L_{photo}$.

$$L_{photo} = \frac{1}{m}\sum((I_{ref} - I_{src}') + (\nabla I_{ref} - \nabla I_{src}')) \cdot M \tag{11}$$

Where $m$ is the sum number of valid points in the mask $M$.

Secondly, the loss of structure similarity (SSIM) $L_{SSIM}$ is set to measure the similarity between $I_{ref}$ and $I_{src}'$. The operation $S$ will be 1 when $I_{ref}$ is the same as $I_{src}'$.

$$L_{SSIM} = \frac{1}{m}\sum \frac{1 - S(I_{ref}, I_{src}')}{2} \cdot M \tag{12}$$

Thirdly, the smooth of final refined depth map can make it less steep in the first-order domain and the second-order domain.

$$L_{smooth} = \frac{1}{n}\sum(e^{-\alpha_2|\nabla I_{ref}|}\left|\nabla \widetilde{Z}_i\right| + e^{-\alpha_3|\nabla^2 I_{ref}|}\left|\nabla^2 \widetilde{Z}_i\right|) \tag{13}$$

Where $n$ is the sum number of points in reference image $I_{ref}$.

Finally, the total pixel-wise loss $L_{pixel}$ can be illustrated as below:

$$L_{pixel} = \lambda_1 L_{photo} + \lambda_2 L_{SSIM} + \lambda_3 L_{smooth} \tag{14}$$
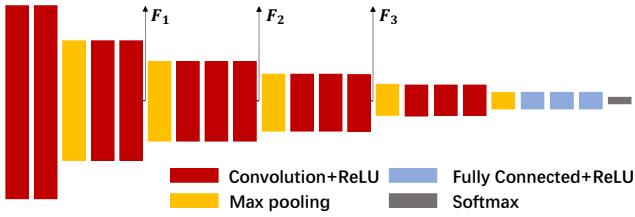
**Figure 5: Feature-wise extraction from pre-trained VGG16**

*3.3.2 Feature-wise Loss.* The network only using pixel-wise loss performs mismatch errors in some challenging scenarios such as textureless and texture repeat areas. In addition to the pixel-wise loss, one of the main improvements of M³VSNet is the use of feature-wise loss. Just like in image style transfer, perceptual loss combining with per-pixel loss improves the performance of style transfer in quality [16]. The feature-wise loss will utilize more semantic information for matching correspondences.

Due to the strong correlation between the estimated depth and pyramid feature network mentioned in section 3.1.1, the high-level feature is extracted from pre-trained VGG16 instead of the pyramid feature network. Through the pre-trained VGG16 network, shown in figure 5, the reference image $I_{ref}$ can extract more semantic high-level information to construct the feature-wise loss function. Here, we extract the layer 8, 15 and 22, which are one half, a quarter and one-eighth the size of the original input images. As a matter of fact, layer 3 output the same size of the original input image, which is actually the reuse of pixel-wise loss function.

For every feature from the VGG16, we construct the loss based on the concept of crossing multi-views. Being similar to section 3.3.1, the corresponding pixel $p_i'$ in $F_{src}$ can be available. The matching features from $F_{ref}$ to $F_{src}$ can be presented as below:

$$F_{src}' = F_{src}(p_i') \qquad (15)$$

The feature domain has a bigger receptive field, which is inspired by that human visual system perceives the scene by its features rather than a single pixel. Therefore, the obstacle of non-ideal areas can be relieved to some extent. The estimated final depth will detect the similarity of features beyond pixel texture value, which benefits from semantic information. The loss $L_F$ is represented as below:

$$L_F = \frac{1}{m} \sum (F_{ref} - F_{src}') \cdot M \qquad (16)$$

The final feature-wise loss function is a weighted sum of different scale of features, which raises the robustness and completeness of point cloud reconstruction. $L_{F_8}$ represents the feature of layer 8 from pre-trained VGG16.

$$L_{feature} = \beta_1 L_{F_8} + \beta_2 L_{F_{15}} + \beta_3 L_{F_{22}} \qquad (17)$$

## 4 EXPERIMENTS

We conduct abundant experiments of our proposed M³VSNet on different datasets. Firstly, we evaluate M³VSNet on the *DTU* dataset and our method outperforms all the previous unsupervised MVS network [7, 17]. Then the ablation studies are carried out to find out potential improvements from our proposed different modules in section 4.3. At last, we test M³VSNet on the *Tanks and Temples* benchmark to verify the generalization ability of our model.

### 4.1 Performance on *DTU*

The DTU dataset is a multi-view stereo dataset that has 124 different scenes with 49 scans for each scene, which is collected by the robotic arms [14][1]. With the lighting change, each scan has seven conditions with the known pose. We use the same train-validation-test split as in MVSNet [31] and MVS² [7]. That is to say, the scenes 1, 4, 9, 10, 11, 12, 13, 15, 23, 24, 29, 32, 33, 34, 48, 49, 62, 75, 77, 110, 114, 118 are selected as the test lists.

*4.1.1 Implementation Detail.* M³VSNet is implemented by Pytorch [? ]. During the training phase, we only use the *DTU*'s training set without any ground-truth depth maps. The resolution of the input image is the crop version of the original picture. That is 640 × 512. Due to the pyramid feature aggregation, the resolution of the final depth is 160 × 128. Additionally, the hypothetical range of depth is sampled from 425mm to 935mm and the depth sample number $D$ is set to 192. The model is trained with the batchsize as 4 in four NVIDIA RTX 2080Ti. By the pattern of data-parallel, each GPU with around 11G available memory could deal with the multi-batch. By using adam optimizer for 10 epochs, the learning rates are set to 1e-3 for the first epoch and decrease by 0.5 for every two epochs. For the balance of different weights in loss, we set $\gamma_1 = 1$, $\gamma_2 = 1$, $\alpha_1 = 0.1$, $\alpha_2 = 0.5$, $\alpha_3 = 0.5$, $\lambda_1 = 0.8$, $\lambda_2 = 0.2$, $\lambda_3 = 0.067$. Beyond that, $\beta_1 = 0.2$, $\beta_2 = 0.8$, $\beta_3 = 0.4$. During each iteration, one reference image and two source images are used. During the testing phase, the resolution of input image is 1600 × 1200.

*4.1.2 Results on DTU.* The official metrics [14] are used to evaluate M³VSNet' performance on the *DTU* dataset. There are three metrics called accuracy, completeness and overall. The overall is the mean value of accuracy and completeness. To prove the effectiveness of M³VSNet, we compare M³VSNet with three classic traditional methods such as Furu [10], Tola [26] and Colmap [23], and with two classic supervised learning-based methods such as SurfaceNet [15] and MVSNet [31], and with the other two existed unsupervised learning-based methods such as Unsup_MVS [17] and MVS² [7].

As shown in the table 1, our proposed M³VSNet outperforms the existed two unsupervised learning-based methods [7, 17]. M³VSNet surpasses Unsup_MVS [17] in all metrics and surpasses MVS² in accuracy and overall except completeness. Therefore, our proposed M³VSNet establishes the state-of-the-arts unsupervised learning methods for multi-view stereo reconstruction. Moreover, M³VSNet surpasses the supervised learning-based MVSNet [31] with the same setting depth hypothesis $D = 192$ in terms of the overall performance of point cloud reconstruction. Compared with traditional MVS methods [9, 23, 26], our proposed M³VSNet achieves significant improvement on the completeness of point cloud reconstruction and outperforms Furu [10] and Tola [26] on the overall quality except Colmap [23] but with high efficiency. For more detailed information in point cloud reconstruction, figure 6 illustrates the qualitative comparison. The reconstruction by M³VSNet has more complete texture details than that without feature-wise loss. With the aid of multi-metric, M³VSNet is more robust so that it
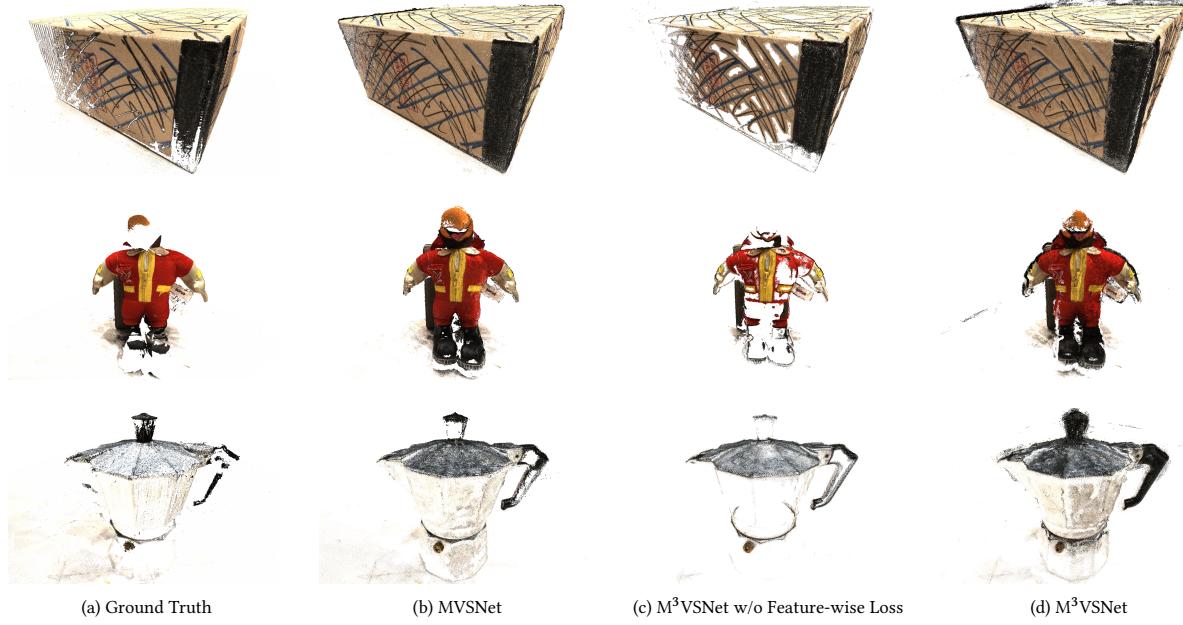
(a) Ground Truth          (b) MVSNet          (c) M$^3$VSNet w/o Feature-wise Loss          (d) M$^3$VSNet

**Figure 6: Qualitative comparison of 3D reconstruction between M$^3$VSNet and supervised methods on the *DTU* dataset. From left to right: ground truth, MVSNet [31], M$^3$VSNet without feature-wise loss and M$^3$VSNet. Our proposed M$^3$VSNet establishes the state-of-the-arts unsupervised learning-based method and achieves comparable performance with MVSNet [31].**

**Table 1: Quantitative results on the *DTU*âĂŹs evaluation set. Three classical MVS methods, two supervised learning-based MVS methods and three unsupervised methods using the distance metric (lower is better) are listed.**

| Method | Mean Distance (mm) | | |
|---|---|---|---|
| | Acc. | Comp. | overall. |
| Furu [10] | 0.612 | 0.939 | 0.775 |
| Tola [26] | **0.343** | 1.190 | 0.766 |
| Colmap [23] | 0.400 | **0.664** | **0.532** |
| SurfaceNet [15] | 0.450 | 1.043 | 0.746 |
| MVSNet(D=192) | **0.444** | **0.741** | **0.592** |
| Unsup_MVS [17] | 0.881 | 1.073 | 0.977 |
| MVS$^2$ [7] | 0.760 | **0.515** | 0.637 |
| **M$^3$VSNet(D=192)** | **0.636** | 0.531 | **0.583** |

recovers more textureless or texture repeat areas and achieves comparable visual performance with original MVSNet [31].

## 4.2 Comparison With Unsupervised Methods

M$^3$VSNet establishes the state-of-the-art unsupervised learning-based MVS network by outperforming the other two existing unsupervised MVS networks [7, 17]. One is unsup_mvs [17], which is almost the first try in this direction but with poor performance where the overall mean distance is 0.977. The other one is MVS$^2$ [7]. Although MVS$^2$ can get a little bit better completeness than

M$^3$VSNet and can reach to 0.637 in the overall mean distance, it consumes more GPU memory due to three cost volumes and regularization needed to be constructed, which is unaffordable for a single NVIDIA RTX 2080Ti used in M$^3$VSNet. As a result, our proposed unsupervised method achieves the best performance on the overall quality of point cloud reconstruction with high efficiency where the accuracy of point cloud is significantly improved.

## 4.3 Ablation Studies

The section begins to analyze the effect of different modules proposed in M$^3$VSNet. There are mainly three contrast experiments carried out. We will explore the effect of pyramid feature aggregation, normal-depth consistency and multi-metric loss.

*Pyramid feature aggregation*. The module can catch more contextual information from low-level to high-level representations. We use the feature pyramid aggregation to output the 1/4 feature. By pyramid feature aggregation, the matching correspondences will be guaranteed to a large extent. As shown in table 2, this module will improve the metric of accuracy and completeness in mean distance. Further, pyramid feature aggregation improves 2% in overall.

*Normal-depth consistency*. Based on the orthogonality between local surface tangent and normal, normal-depth consistency is introduced to regularize the depth in 3D space. Absolute depth error is used to evaluate the quality of estimated depth. Here we use the percentage of depth error within 2mm, 4mm, and 8mm compared with ground-truth depth maps (Higher is better). As shown in table 3, the performance with the aid of normal-depth consistency surpasses that without normal-depth consistency in all metrics.

**Table 2: Comparison of the performance in pyramid feature aggregation using the distance metric (lower is better).**

| Method | Mean Distance (mm) | | |
|---|---|---|---|
| | Acc. | Comp. | overall |
| without pyramid feature aggregation | 0.638 | 0.554 | 0.596 |
| **with pyramid feature aggregation** | **0.636** | **0.531** | **0.583** |

**Table 3: Comparison of the performance in normal-depth consistency using the depth error (higher is better).**

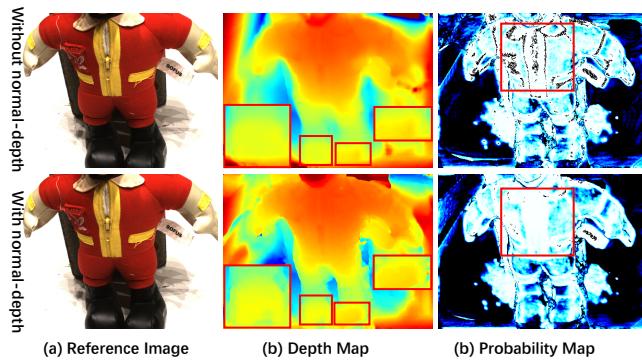| Depth Error (mm) | % < 2 | % < 4 | % < 8 |
|---|---|---|---|
| without normal-depth consistency | 58.8 | 74.8 | 83.8 |
| **with normal-depth consistency** | **60.3** | **76.9** | **85.7** |



**Figure 7: Qualitative comparison of normal-depth consistency in depth maps and probability maps (White is 100%, black is 0%). The performance of M³VSNet will be more robust and more accurate even in non-ideal environments.**

Figure 7 demonstrates the comparison with and without normal-depth consistency. For the same reference image, the depth map with normal-depth consistency is more accurate than that without normal-depth consistency. The depth of foot (the part framed in red) in the first row is more ambiguous than that in the second row along the edge and the plane. What's more, the foot part is textureless with reflective effect. Therefore, the robustness of performance with normal-depth consistency will be guaranteed even in non-ideal environments. Normal-depth consistency will make the estimated depth more precise and reasonable in 3D space. Furthermore, in probability maps, the probability of small local areas such as the collar and the zipper (the part framed in red) will be high with the aid of normal-depth consistency. The module also improves the probability of correct matching correspondences in some different planes. Figure 7 and table 3 can prove the significant benefits of normal-depth consistency for accuracy and continuity.

***Multi-metric loss.*** Multi-metric loss contains pixel-wise loss and feature-wise loss, which learns the inherent constraints from different perspectives of matching correspondences. The try to

**Table 4: Comparison of the performance in the different loss terms using the percentage of depth error (higher is better).**

| Depth Error (mm) | % < 2 | % < 4 | % < 8 |
|---|---|---|---|
| B | 22.1 | 36.5 | 50.8 |
| B+G | 25.2 | 40.7 | 55.3 |
| B+G+SSIM | 27.5 | 44.2 | 58.8 |
| B+G+SSIM+Smooth | 57.5 | 75.2 | 85.4 |
| **Multi-metric loss** | **60.3** | **76.9** | **85.7** |

**Table 5: Comparison of the performance in different loss (lower is better). The scale of 1/2 represents that the feature (corresponding to layer 8) extracted from the pre-trained VGG16 networks is half the size of the original reference image. The scales of 1/4, 1/8, 1/16 correspond to layer 15, 22, 29.**

| Method | Mean Distance (mm) | | |
|---|---|---|---|
| | Acc. | Comp. | overall |
| only pixel-wise | 0.832 | 0.924 | 0.878 |
| pixel-wise+1/4 feature | 0.646 | 0.591 | 0.618 |
| **pixel-wise+1/2,1/4,1/8 feature** | **0.636** | **0.531** | **0.583** |
| pixel-wise+1/2,1/4,1/8,1/16 feature | 0.566 | 0.653 | 0.609 |

feature-wise loss is effective in previous related works [16][29][4]. We have compared the performance of different combinations of pixel-wise loss and feature-wise loss. What's more, how to select the multi-scale features is also taken into consideration.

In the ablation study of loss terms in pixel-wise loss, the absolute depth error is used to evaluate the performance of different loss terms. As shown in table 4, B is the baseline with the only photometric loss in pixel level. G represents the gradient consistency loss. As demonstrated in section 3.3, SSIM is represented as $L_{SSIM}$ and smooth is represented as $L_{smooth}$. The terms of G and SSIM improve the results slightly and the term of Smooth contributes a lot with effective improvement. In general, it's apparent that the proposed each loss will improve the performance of M³VSNet.

In the ablation study of loss terms in feature-wise loss, as illustrated in table 5, the overall of only pixel-wise loss is relatively higher (lower is better). Besides, the different combinations of feature-wise losses make it an impressive improvement. We do some ablation studies on the different combinations of features from pre-trained VGG16. Adding the 1/16 feature improves the accuracy but deteriorate the completeness. By comparison, the combination of 1/2, 1/4, 1/8 features achieves the best result.

## 4.4 Generalization Ability on *Tanks & Temples*

To evaluate the generalization ability of our proposed M³VSNet, we use the intermediate *Tanks and Temples* benchmark that has high-resolution images of outdoor large-scale scenes. The model of our proposed M³VSNet trained on the *DTU* dataset is transferred to the *Tanks & Temples* benchmark without any finetuning. The intermediate *Tanks and Temples* benchmark contains kinds of images with the resolution of 1920 × 1056 and with the depth

**Table 6: Quantitative comparison of point cloud reconstruction on the *Tanks and Temples* benchmark (higher is better). M³VSNet surpasses other unsupervised methods by the mean score in the leaderboard of intermediate *T&T* [18].**

| Method | Mean | Family | Francis | Horse | Lightouse | M60 | Panther | Playground | Train |
|---|---|---|---|---|---|---|---|---|---|
| **M³VSNet** | **37.67** | **47.74** | **24.38** | 18.74 | 44.42 | 43.45 | 44.95 | **47.39** | **30.31** |
| MVS² | 37.21 | 47.74 | 21.55 | **19.50** | **44.54** | **44.86** | **46.32** | 43.48 | 29.72 |



Family  Francis  Horse  M60



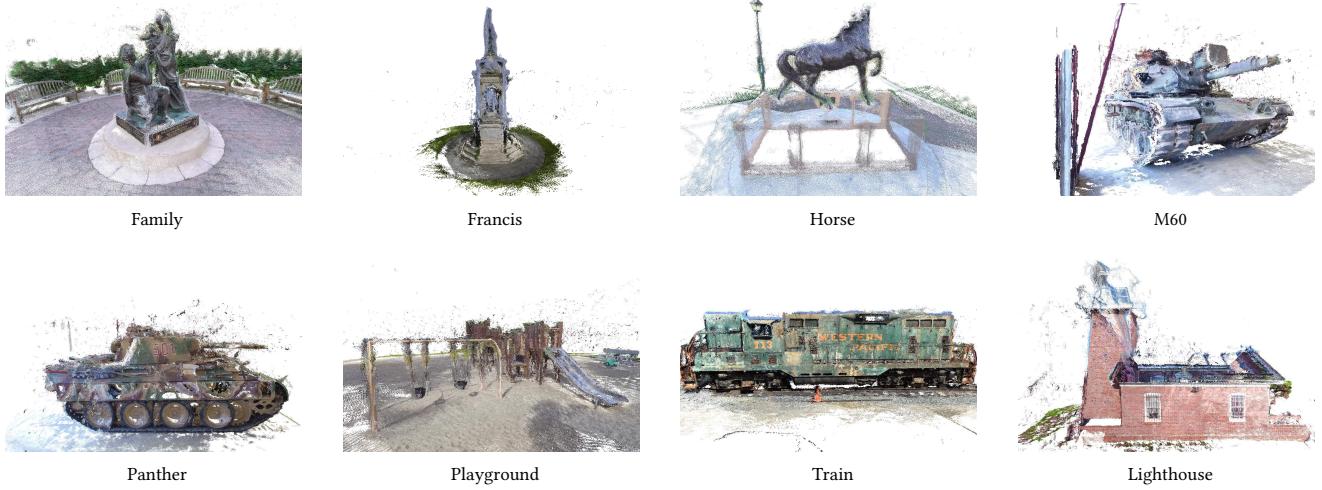Panther  Playground  Train  Lighthouse

**Figure 8: The performance of M³VSNet on the *Tanks and Temples* benchmark [18] without any finetuning. The quality of dense point cloud reconstruction in large-scale scene shows the powerful generalization ability of M³VSNet.**

hypothesis $D = 160$. Another core hyperparameter is the photometric threshold in the process of depth fusion. For the same depth maps, the different photometric thresholds will lead to different performances. Higher photometric threshold will cause better accuracy but worse completeness. In turn, lower photometric threshold will introduce better completeness but worse accuracy. For our proposed M³VSNet, the photometric threshold is set to 0.6 and we get the following results. As shown in table 6, the ranking is selected from the leaderboard of the intermediate *Tanks and Temples* benchmark. Our proposed M³VSNet is better than MVS² by the mean score of 8 scenes, which is the best unsupervised MVS network until April 17, 2020. In table 6, the higher the mean score, the higher the ranking relatively. Further, the score of M³VSNet in Playground scene is 47.39, which is better than the score of MVS² 43.48. Therefore, our proposed M³VSNet ranks higher than MVS². As a matter of fact, the final ranking is the mean of the independent ranking of 8 scenes, which is different from the calculation method of the mean score. The dense point clouds are presented in figure 8, which are reasonable and complete for Family, Francis, Horse, M60, Panther, Playground, Train, Lighthouse scenes. Besides, the robustness of our proposed M³VSNet also play an important role in non-ideal areas in *Tanks and Temples* benchmark such as the sand in the scene of Playground. In view of the above, the performance in table 6 and figure 8 demonstrates the powerful generalization ability of our proposed M³VSNet.

## 5 CONCLUSION

In this paper, we propose an unsupervised multi-metric network for multi-view stereo reconstruction named M³VSNet, which improve the robustness and completeness of point cloud even in non-ideal environments. The proposed novel multi-metric loss function, namely pixel-wise and feature-wise loss function, can capture more semantic information to learn the inherent constraints from different perspectives of matching correspondences. The performance of point cloud reconstruction in non-ideal environments for robustness and completeness will also benefit from the multi-metric loss mainly. Besides, with the incorporation of normal-depth consistency, M³VSNet improves the accuracy and continuity of the estimated depth maps by the orthogonality between normal and local surface tangent. Extensive experiments show that our proposed M³VSNet outperforms the previous state-of-the-arts unsupervised learning-based methods and achieves comparable performance with original MVSNet [31] on the *DTU* dataset and demonstrates the powerful generalization ability on the *Tanks & Temples* benchmark [18] with effective improvement. In the future, more MVS datasets with high precision are desired. To relief the high cost of datasets, the domain transfer for different datasets can be improved and enhanced. What's more, multi-task such as object detection, semantic and instance segmentation, depth completion, etc. can be combined with multi-view stereo reconstruction for the time to come.

# REFERENCES

[1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. 2016. Large-Scale Data for Multiple-View Stereopsis. *International Journal of Computer Vision* 120 (2016), 153–168.

[2] Ibraheem Alhashim and Peter Wonka. 2018. High quality monocular depth estimation via transfer learning. *arXiv preprint arXiv:1812.11941* (2018).

[3] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B Goldman. 2009. PatchMatch: A randomized correspondence algorithm for structural image editing. In *ACM Transactions on Graphics (ToG)*. ACM, 24.

[4] Wang Benzhang, Feng Yiliu, Fu Huini, and Hengzhu Liu. 2018. Unsupervised Stereo Depth Estimation Refined by Perceptual Loss. In *2018 Ubiquitous Positioning, Indoor Navigation and Location-Based Services (UPINLBS)*. IEEE, 1–6.

[5] Neill DF Campbell, George Vogiatzis, Carlos Hernández, and Roberto Cipolla. 2008. Using multiple hypotheses to improve depth-maps for multi-view stereo. In *European Conference on Computer Vision*. Springer, 766–779.

[6] Hainan Cui, Xiang Gao, Shuhan Shen, and Zhanyi Hu. 2017. HSfM: Hybrid structure-from-motion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1212–1221.

[7] Yuchao Dai, Zhidong Zhu, Zhibo Rao, and Bo Li. 2019. MVS2: Deep Unsupervised Multi-View Stereo with Multi-View Symmetry. In *2019 International Conference on 3D Vision (3DV)*. IEEE, 1–8.

[8] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. 2015. Flownet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*. 2758–2766.

[9] Jan-Michael Frahm, Pierre Fite-Georgel, David Gallup, Tim Johnson, Rahul Raguram, Changchang Wu, Yi-Hung Jen, Enrique Dunn, Brian Clipp, Svetlana Lazebnik, et al. 2010. Building rome on a cloudless day. In *European Conference on Computer Vision*. Springer, 368–381.

[10] Yasutaka Furukawa and Jean Ponce. 2007. Accurate, Dense, and Robust Multi-View Stereopsis. *2007 IEEE Conference on Computer Vision and Pattern Recognition* (2007), 1–8.

[11] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. 2015. Massively parallel multiview stereopsis by surface normal diffusion. In *Proceedings of the IEEE International Conference on Computer Vision*. 873–881.

[12] Michael Goesele, Brian Curless, and Steven M Seitz. 2006. Multi-view stereo revisited. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, Vol. 2. IEEE, 2402–2409.

[13] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. 2019. Cascade Cost Volume for High-Resolution Multi-View Stereo and Stereo Matching. *arXiv preprint arXiv:1912.06378* (2019).

[14] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. 2014. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 406–413.

[15] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. 2017. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *Proceedings of the IEEE International Conference on Computer Vision*. 2307–2315.

[16] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*. Springer, 694–711.

[17] Tejas Khot, Shubham Agrawal, Shubham Tulsiani, Christoph Mertz, Simon Lucey, and Martial Hebert. 2019. Learning Unsupervised Multi-View Stereopsis via Robust Photometric Consistency. *arXiv preprint arXiv:1905.02706* (2019).

[18] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. 2017. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1–13.

[19] Hamid Laga. 2019. A Survey on Deep Learning Architectures for Image-based Depth Reconstruction. *arXiv preprint arXiv:1906.06113* (2019).

[20] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.

[21] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. 2019. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *Proceedings of the IEEE International Conference on Computer Vision*. 10452–10461.

[22] Reza Mahjourian, Martin Wicke, and Anelia Angelova. 2018. Unsupervised learning of depth and ego-motion from monocular video using 3d geometric constraints. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5667–5675.

[23] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. 2016. Pixelwise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision*. Springer, 501–518.

[24] Steven M Seitz, Brian Curless, James Diebel, Daniel Scharstein, and Richard Szeliski. 2006. A comparison and evaluation of multi-view stereo reconstruction algorithms. In *2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06)*, Vol. 1. IEEE, 519–528.

[25] Sudipta N Sinha, Philippos Mordohai, and Marc Pollefeys. 2007. Multi-view stereo via graph cuts on the dual of an adaptive tetrahedral mesh. In *2007 IEEE 11th International Conference on Computer Vision*. IEEE, 1–8.

[26] Engin Tola, Christoph Strecha, and Pascal Fua. 2011. Efficient large-scale multi-view stereo for ultra high-resolution image sets. *Machine Vision and Applications* 23 (2011), 903–920.

[27] Benjamin Ummenhofer, Huizhong Zhou, Jonas Uhrig, Nikolaus Mayer, Eddy Ilg, Alexey Dosovitskiy, and Thomas Brox. 2017. Demon: Depth and motion network for learning monocular stereo. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5038–5047.

[28] Hoang-Hiep Vu, Patrick Labatut, Jean-Philippe Pons, and Renaud Keriven. 2011. High accuracy and visibility-consistent dense multiview stereo. *IEEE transactions on pattern analysis and machine intelligence* 34, 5 (2011), 889–901.

[29] Anjie Wang, Zhijun Fang, Yongbin Gao, Xiaoyan Jiang, and Siwei Ma. 2018. Depth estimation of video sequences with perceptual losses. *IEEE Access* 6 (2018), 30536–30546.

[30] Zhenheng Yang, Peng Wang, Wei Xu, Liang Zhao, and Ramakant Nevatia. 2018. Unsupervised learning of geometry from videos with edge-aware depth-normal consistency. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

[31] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. 2018. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 767–783.

[32] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. 2019. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5525–5534.

[33] Hongwei Yi, Zizhuang Wei, Mingyu Ding, Runze Zhang, Yisong Chen, Guoping Wang, and Yu-Wing Tai. 2019. Pyramid Multi-view Stereo Net with Self-adaptive View Aggregation. *arXiv preprint arXiv:1912.03001* (2019).

[34] Tinghui Zhou, Matthew Brown, Noah Snavely, and David G Lowe. 2017. Unsupervised learning of depth and ego-motion from video. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1851–1858.