

# MVPSNet: Fast Generalizable Multi-view Photometric Stereo

Dongxu Zhao<sup>1</sup> Daniel Lichy<sup>2</sup> Pierre-Nicolas Perrin<sup>1</sup> Jan-Michael Frahm<sup>1</sup> Soumyadip Sengupta<sup>1</sup>

<sup>1</sup>University of North Carolina at Chapel Hill <sup>2</sup>University of Maryland, College Park

{dongxuz1, pnperrin, jmf, ronisen}@cs.unc.edu dlichy@umd.edu

## Abstract

We propose a fast and generalizable solution to Multi-view Photometric Stereo (MVPS), called MVPSNet. The key to our approach is a feature extraction network that effectively combines images from the same view captured under multiple lighting conditions to extract geometric features from shading cues for stereo matching. We demonstrate these features, termed ‘Light Aggregated Feature Maps’ (LAFM), are effective for feature matching even in textureless regions, where traditional multi-view stereo methods fail. Our method produces similar reconstruction results to PS-NeRF, a state-of-the-art MVPS method that optimizes a neural network per-scene, while being  $411 \times$  faster (105 seconds vs. 12 hours) in inference. Additionally, we introduce a new synthetic dataset for MVPS, sMVPS, which is shown to be effective to train a generalizable MVPS method.

## 1. Introduction

3D reconstruction of an object can be achieved either through camera viewpoint variations, Multi-view Stereo (MVS), or by lighting direction variations, Photometric Stereo (PS). Both MVS and PS have relative strengths and weaknesses. While MVS succeeds in obtaining accurate global shapes, it suffers in textureless regions due to poor feature matching, often resulting in reconstructions that lack local details. On the other hand, PS produces accurate local details, even in textureless regions, by using shading information but fails to reconstruct accurate global shapes. In this paper, we focus on the problem of Multi-view Photometric Stereo (MVPS) where both camera viewpoint and lighting direction variations are used to accurately reconstruct global and local details of a 3D shape, even in textureless regions.

3D reconstruction techniques that produce high-quality results using only viewpoint variations (MVS) rely on test-time optimization, often by training neural networks per scene [57, 68, 71]. These methods are computationally inefficient, typically taking hours of compute time on a high-end GPU for each object. Existing MVS methods

[22, 60, 66] that focus on computational efficiency employ feed-forward neural networks that are efficient but fail to produce high-quality details, especially in textureless regions. Existing MVPS approaches can produce high-quality reconstructions but require computationally inefficient per-scene training or optimization [30, 31, 32, 65]. Sometimes additional manual efforts and carefully crafted refinement steps are also needed [38, 51]. In contrast, we propose an efficient feed-forward neural architecture, MVPSNet, that can generalize to unseen objects and achieve similar reconstruction quality to that of per-scene optimization techniques while being computationally efficient during inference.

We design MVPSNet by taking inspiration from various deep MVS architectures [8, 15, 20, 22, 66] that are generalizable, computationally efficient, and can operate on high-resolution images. However, these approaches fail in textureless regions, and their reconstructed meshes often lack details. We choose the CasMVSNet [22] architecture as our feature matching module, which has been repeatedly used by various MVS pipelines [8, 15, 20] for its simplicity and efficiency, and augment it to effectively incorporate lighting variation cues for better prediction of 3D shapes. To our knowledge, we are the first to propose a feed-forward generalizable approach to Multi-view Photometric Stereo.

We introduce a multi-scale feature representation, called Light Aggregated Feature Maps (LAFM), whose role is to extract detailed geometric features from images by utilizing lighting variations. For brevity, we define **Multi-light Images** as a collection of images taken from the same viewpoint under different directional lighting conditions. Our intuition is that LAFM can efficiently aggregate shading patterns from multi-light images, by creating an ‘artificial shading texture’ in the textureless region. Multiscale LAFMs will then be used to construct a sequence of Cost Volumes to match features across sparse viewpoints in order to predict a depth map for each viewpoint. We also predict surface normals from LAFM for each viewpoint, enabling us to capture features related to high-frequency local details.

We further show that the surface normal predicted by using LAFM can be used in addition to the depth maps to produce a more detailed mesh than using the depth maps alone.

To train the proposed MVPSNet architecture, we introduce a new synthetic MVPS dataset. Our synthetic dataset consists of shapes from sculpture dataset [61] and random compositions of primitive shapes generated by [64]. We render these shapes with spatially varying Cook-Torrance BRDF under different camera viewpoints and lighting directions. We train MVPSNet on these rendered images with ground-truth supervision over predicted depth and surface normal maps. The trained model generalizes to real-world test scenes from DiLiGenT-MV [38] dataset. We show that simply re-training CasMVSNet on our dataset improves reconstruction quality over the pre-trained model on DiLiGenT-MV by 32%, proving the effectiveness of our synthetic MVPS dataset for generalization.

We evaluate our approach on the only publicly available MVPS benchmark, the DiLiGenT-MV [38] dataset. Compared to the state-of-the-art MVPS technique, PS-NeRF [65], which optimizes a neural network per-scene, our proposed MVPSNet is  $\sim 411 \times$  faster (105 seconds vs 12 hours) while producing similar reconstruction quality (L1 Chamfer distance of 0.82 vs 0.81, F-score on L2 distance of 0.985 vs 0.983). We further show that adding LAFM features significantly improves reconstruction quality over CasMVSNet by 34% in L1 Chamfer distance. We also observe that refining the reconstructed mesh derived from depth maps with predicted surface normals from LAFM features improves reconstruction quality as shown in Fig 3.

In summary, the key contributions of this paper include:

- Light Aggregated Feature Maps (LAFM) that can efficiently utilize multi-light images to extract detailed geometric features, especially in textureless regions. The surface normal predicted from LAFM also improves mesh reconstruction quality.
- A synthetic MVPS dataset for training generalizable MVPS methods, which also improves CasMVSNet by 32%.
- A fast and generalizable Multi-view Photometric Stereo pipeline that is  $411 \times$  faster while producing similar reconstruction accuracy compared to state-of-the-art per-scene optimization approach [65].

## 2. Related work

**Multi-View Stereo (MVS).** MVS is a 3D reconstruction technique that utilizes multiple images captured from different viewpoints. While various techniques for MVS have been proposed, one commonly used approach that is relevant to our work involves constructing cost volumes similar to Plane Sweeping Algorithm [14]. To create a cost volume, features are matched across neighboring viewpoints, and the quality of the features plays a critical role in the final reconstruction quality. Traditional methods [6, 17, 18, 21, 29, 54, 58, 59] use human-defined or hand-crafted image processing operators to extract feature maps.

Method	Generalizable	Mesh Reconstruction
PJ16 [51]	✗	Base mesh+displacement map
LZ20 [38]	✗	3D points+PSR [33]+Optimization [48]
BKW22 [32]	✗	MLP+Marching Cube [44]
BKC22 [30]	✗	MLP+Marching Cube [44]
PS-NeRF [65]	✗	MLP+MISE [46]
BKW23 [31]	✗	MLP+Marching Cube [44]
Ours	✓	3D Points+Screened Poisson [34]

Table 1. Comparison of our method with prior MVPS methods.

With recent advances in deep learning, features are extracted with a deep neural network to build cost volumes and then predict per-view depth maps [23, 28, 45, 63, 66, 67] or disparity map [26].

The most relevant previous works are MVSNet [66] and its variations. MVSNet [66] uses homography to warp feature maps and a 3D CNN to regularize cost volumes. CasMVSNet [23] outperforms MVSNet in terms of accuracy and efficiency by building the 3D cost volume in a cascaded manner. TransMVSNet [16] builds upon CasMVSNet and adopts a transformer to consider intra-image and inter-image feature interactions, which further improves the result of CasMVSNet.

**Photometric Stereo (PS).** PS (introduced in [62]) uses lighting variation to reconstruct 3D shapes from a single viewpoint (see [56] for surveys). Calibrated PS approaches, like Chen *et al.* [11], train a neural network to predict surface normals using data with known lightings. Uncalibrated PS approaches [9, 10, 12] first predict the lighting parameters before solving for surface normals. While most PS works use a large number of images for inference, some use fewer [27, 42], or even one image [7, 39, 40, 55] (often called Shape from Shading). PS approaches are mostly based on feed-forward networks that generalize and can produce near real-time inference with low computational cost [41].

**Multi-View Photometric Stereo (MVPS).** MVPS was initially proposed in [25] by combining PS with object silhouettes to reconstruct textureless shiny objects with fine details. However, this method only works well for specific parametric BRDF models [30]. Later, Li *et al.* [38, 73] propose to get iso-depth contours from PS images and sparse 3D points using structure-from-motion, which are propagated to recover complete 3D shape. Park *et al.* [50, 51] use a planar mesh parameterization technique to parameterize a coarse mesh from MVS and take advantage of this 2D parameter domain to perform MVPS. However, these traditional MVPS methods require an initial 3D reconstruction and their performance is sensitive to it. Since they consist of multiple steps, careful execution or expert interventions are often performed to get good results [38, 51].

Recently, inspired by NeRF [47], various algorithms have been proposed that optimize a neural network per-scene for MVPS. Kaya *et al.* [32] train a deep PS network first and condition the color rendering in NeRF [69] on normals predicted from PS. The reconstructed mesh, however, exhibits multiple artifacts. The authors in [30] propose to

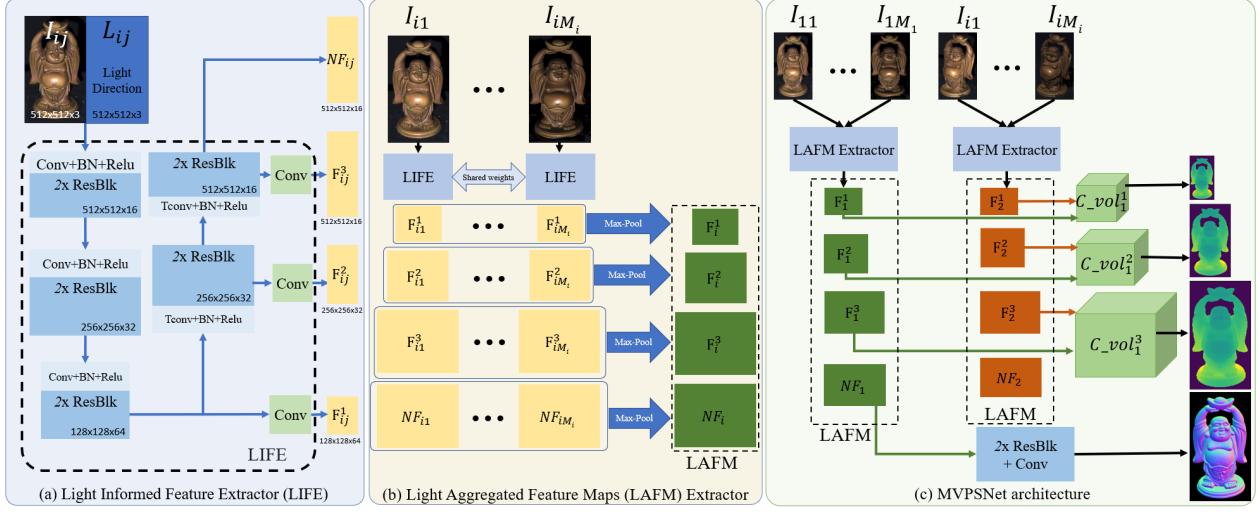


Figure 1. Overview of our network architecture. (a) Light Informed Feature Extractor (LIFE) produces a multi-scale feature representation; (b) Light Aggregated Feature Maps (LAFM) Extractor aggregates these feature across images of varying lighting conditions but same view; and (c) LAFM is used to create Cost Volume and predict depth map, similar to CasMVSNet [23], in addition to normal map.

train a deep PS network and a deep MVS network extended with uncertainty estimation separately, and use these to fit an SDF represented by an MLP. To further enable reconstruction on anisotropic and glossy objects, Kaya *et al.* [31] add a neural volume rendering module to the MLP used in [30] to better fuse PS and MVS measurements. PS-NeRF [65] solves the task of jointly estimating the geometry, materials and lights. It first regularizes the gradient of a UNISURF [49] with estimated normals from PS, and then uses separate MLPs to explicitly model surface normal, BRDF, lights, and visibility which are optimized based on a shadow-aware differentiable rendering layer. Recent works have also used physically based differentiable rendering either inside a NeRF framework [4] or separately for optimization [70]. However, these methods optimize models for each object independently, thus they do not generalize and are computationally inefficient.

In contrast, to our knowledge, we are the first to provide a generalizable solution to Multi-view Photometric Stereo by training on our proposed synthetic MVPS dataset.

### 3. Our approach

#### 3.1. Problem Setup

We focus on the problem of calibrated multi-view photometric stereo, i.e. the locations of the light sources and the cameras are known *a priori* (calibrated prior to capture). The input data consists of a set of multi-light images of an object captured from multiple views.

Concretely, for the  $i$ -th view we have  $M_i$  images with varying lighting directions  $l_{ij}$ , denoted as  $I_{ij}$ . We refer to the collection  $\{I_{i1}, \dots, I_{iM_i}\}$  as the multi-light images for the  $i$ -th view. For each view, we are given camera intrinsic matrix  $K_i$  and camera extrinsic parameters in the form of a

rotation matrix,  $R_i$ , and a translation vector,  $t_i$ . Similar to virtually all MVS methods, we assume that we are provided with the depth range for each view.

#### 3.2. Motivation

Our approach follows a long line of work in Multi-view Stereo which uses Plane Sweep to construct a Cost Volume and predicts a depth map aligned with a reference image. Recent advances in Plane Sweep Stereo using deep neural networks, especially CasMVSNet [23], have proven to be generalizable across scenes and can predict high-resolution reconstruction in a matter of seconds. In contrast to previous Multi-view Photometric Stereo approaches, which optimize a neural network per scene, our goal is to produce a generalizable solution. Thus we aim to build on the Plane Sweep stereo architecture proposed in CasMVSNet [23].

CasMVSNet learns a deep image feature encoder for extracting representative features that can aid in feature matching across multiple views and create a better Cost Volume. However, these features are often ambiguous for non-textured regions and often fail to preserve the geometric details. We believe incorporating lighting variations along with viewpoint variations can lead to better features, which in turn will produce better Cost Volumes and depth maps.

To this end, we introduce ‘Light Aggregated Feature Maps’ (LAFM), whose goal is to extract detailed geometric features, even for textureless regions, by jointly learning to aggregate feature maps over all images captured from a single viewpoint and multiple lighting directions. To obtain effective features that capture geometric details we use LAFM to regress surface normals. We show that LAFM provide superior information for stereo-matching than single-lighting feature maps (as used in CasMVSNet) and thus provide a better reconstruction. We also show that surface normals

predicted using LAFM can be used during mesh reconstruction to improve quality over depth maps alone.

Our approach proceeds in three stages. We first extract multi-scale feature representation from each image, along with its lighting directions, using a shared neural network, ‘Light Informed Feature Extractor’ (LIFE). Then we aggregate features extracted by LIFE across all images captured under the same viewpoint but different lighting conditions using a max-pooling operation to form ‘Light Aggregated Feature Maps’ (LAFM). We can then create a Cost Volume for the reference view by matching LAFM of the reference view with all the LAFM from neighboring views. Finally, for each reference view, we predict a depth map using cost volume regularization and a surface normal map from the LAFM. We train our system in a multi-task learning framework with supervised losses over depth and normal predictions. In the following sections, we provide the details of our MVPSNet pipeline. We provide an overview of MVPSNet architecture in Figure 1.

### 3.3. Light Aggregated Feature Maps (LAFM)

We introduce Light Aggregated Feature Maps (LAFM) that provide geometrically distinct multi-scale features for Cost Volume creation in a Plane Sweep Stereo approach. Our key observation is that the multi-light images provide us with important information for feature matching. For textureless regions, the variation in shading (including cast and attached shadows) created by different lighting directions can be interpreted as ‘artificial’ textures. Thus the role of LAFM for textureless regions is to capture the variation in shading as an ‘artificial’ texture that can be used for feature matching across different viewpoints. We also use LAFM to predict surface normal maps, enabling it to capture geometric details required for producing high-quality normal maps. Hence LAFM can capture better features for textureless regions and for reconstructing details, which were absent in the usual deep image features used in deep multi-view stereo algorithms.

We first define a multi-scale feature extractor, Light Informed Feature Extractor (LIFE), that takes an image  $I_{ij}$  associated with its lighting direction  $L_{ij}$  as input and produces features maps at three different scales  $F_{ij}^1, F_{ij}^2, F_{ij}^3$  at resolutions  $1/4, 1/2, 1$  of the input resolution, and another feature map  $NF_{ij}$  that will be used for normal prediction. The network architecture of LIFE is shown in Fig. 1(a) and will be discussed in details in the supplementary material.

$$NF_{ij}, F_{ij}^1, F_{ij}^2, F_{ij}^3 = LIFE(I_{ij}, l_{ij}; \theta) \quad (1)$$

Note  $L_{ij}$  is of the same resolution as  $I_{ij}$  by simply repeating the same 3-dimensional lighting vector at each pixel.

Then we extract these multi-scale features for every image captured under the same viewpoint and different lighting conditions,  $\{I_{i1}, \dots, I_{iM_i}\}$ , using the same shared encoder LIFE. Let the feature maps obtained from these images be denoted as:  $\{NF_{ij}, F_{ij}^1, F_{ij}^2, F_{ij}^3\}$ ,  $j = 1, \dots, M_i$ .

We create ‘Light Aggregated Feature Maps’ (LAFM) from these multi-scale feature representations by simply performing a max-pooling operation for each scale, as proposed in [9], as:

$$F_i^s = \max_j F_{ij}^s, \quad \forall s = 1, 2, 3 \quad (2)$$

$$NF_i = \max_j NF_{ij}. \quad (3)$$

Thus for multi-light images we obtain LAFM as  $LF_i = \{NF_i, F_i^1, F_i^2, F_i^3\}$ .

The features at 3 scales  $F_i^1, F_i^2, F_i^3$  are then used to build cost volumes using differentiable homography warping, which we will talk about in detail in Section 4.1. The normal feature  $NF_i$  is fed into a lightweight normal regression network to predict per-view normal map, as shown in Figure 1(c). With the supervision from normal information and depth information, our LAFM benefit from the advantages of both MVS and PS which are good at global shape modeling and high-frequency component reconstruction, respectively.

### 3.4. Cost Volume and Depth Map Prediction

Given Light Aggregated Feature Maps (LAFM),  $LF_i$ , for each view  $i$ , we aim to build a cost volume for each reference view by selecting a set of source views with sufficient overlap. We adopt the multi-scale cost volume construction proposed in CasMVSNet [23], where the plane sweep is first performed at a low resolution and then at higher resolutions. Depth estimated from the previous step is used for generating depth proposals for the next step. Multi-scale cost volume reconstruction and depth map prediction follow the following steps.

**Step 1: Depth hypothesis generation.** We generate hypothesis depths for each pixel based on the lower resolution depth estimated at the previous resolution. We store these in  $h$ , where

$$h(u, v, w) = Up(D^{s-1})(u, v) + \Delta_s \left( \frac{w}{N_s - 1} - \frac{1}{2} \right). \quad (4)$$

Here  $h(u, v, w)$  is the  $w$ -th depth hypothesis at pixel  $(u, v)$ .  $Up(D^{s-1})$  is the depth map at the previous lower resolution upsampled to the current resolution.  $\Delta_s$  is the length of the depth interval we are searching at scale  $s$ .  $N_s$  is the number of hypothesis depths at the current scale.

**Step 2: Cost-Volume construction.** The cost volume is a way of robustly searching for matches between a point  $u, v$  on a reference image  $I_r$  and a point on the corresponding epipolar line in the source images  $I_{s_k}$ . Concretely, consider a pixel  $(u, v)$  in the reference image. For every hypothesis depth  $d$ , we get a corresponding point in the source image on the epipolar line for  $(u, v)$ . We denote this point by  $(u', v') = \text{warp}_{rs_k}(u, v, d)$  where

$$\begin{bmatrix} u' \\ v' \\ 1 \end{bmatrix} \sim K_{s_k} R_{s_k}^T \left[ \left( R_r K_r^{-1} d \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} + t_r \right) - t_s \right]. \quad (5)$$

We then construct a per-image volume:

$$F_{\text{vol}}^s(u, v, w) = F_i^s(\text{warp}(u, v, h(u, v, w))), \quad (6)$$

where  $i$  runs over the reference and source views i.e.  $i \in \{r, s_1, \dots, s_k\}$ . These volumes are then aggregated into a single cost volume by taking their variance, which checks for the photo-consistency of the depth proposal  $d$  for pixel  $(u, v)$  in the reference image and the corresponding pixels in the warped sources images  $s_k$ :

$$\text{agg\_vol}^s(u, v, w) = \text{var}_i(F_{\text{vol}}(u, v, w)_i^s). \quad (7)$$

**Step 3: Cost-Regularization.** In this step we pass the aggregated volume through a 3D convolutional network and take a softmax to convert it to a match probability, using

$$\text{prob\_vol} = \text{soft\_max}_w(\text{reg\_net}^s(\text{agg\_vol}^s)). \quad (8)$$

**Step 4: Regression.** We take the expectation of the hypothesis depths over the match probability given by the probability volume to obtain the depth at the current scale:

$$D^s(u, v) = \sum_w \text{prob\_vol}(u, v, w) h(u, v, w). \quad (9)$$

This whole process is summarized in algorithm 1

#### Algorithm 1 MVPSNet Algorithm

---

```

1:  $NF_{ij}, F_{ij}^1, F_{ij}^2, F_{ij}^3 = \text{LIFE}(I_{ij}, l_{ij}; \theta)$ 
2:  $NF_i = \max_j NF_{ij}; F_i^s = \max_j F_{ij}^s \quad \forall s = 1, 2, 3.$ 
3:  $N_i = \text{normal\_regression\_net}(NF_i)$ 
4:  $D^0(u, v) = (\text{max\_depth} + \text{min\_depth})/2$ 
5: for  $s = 1$  to  $3$  do
6:    $h(u, v, w) = Up(D^{s-1}(u, v)) + \Delta_s(\frac{w}{N_w - 1} - \frac{1}{2})$ 
7:    $F_{\text{vol}}^s(u, v, w) = F_i^s(\text{warp}(u, v, h(u, v, w)))$ 
8:    $\text{agg\_vol}^s(u, v, w) = \text{var}_i(F_{\text{vol}}_i^s)$ 
9:    $\text{prob\_vol} = \text{soft\_max}_w(\text{reg\_net}^s(\text{agg\_vol}^s))$ 
10:   $D^s(u, v) = \sum_w \text{prob\_vol}(u, v, w) h(u, v, w)$ 
11: end for
```

---

Once we have depth and surface normal for each view, our mesh reconstruction pipeline consists of three steps: depth filtering, lifting depth and normal maps to point cloud, and reconstructing mesh using Screened Poisson [34] (See supplementary materials for details).

### 3.5. Synthetic MVPS dataset

A key component of our method is that we can learn better features for stereo matching, especially in textureless regions, by learning features that incorporate multi-lighting cues. However, there is no existing MVPS dataset that is large enough for neural network training. Therefore, we generate a large-scale synthetic dataset, sMVPS, consisting of two sub-datasets sMVPS-sculpture (800 train scenes/4 test scenes) and sMVPS-random (1000 train scenes/20 test scenes).

sMVPS-sculpture consists of objects from the sculpture dataset [61] while sMVPS-random includes objects composed of random primitives from [64]. The objects were generated following the method of [43] with spatially varying Cook-Torrance BRDF. We render images from 20 viewpoints surrounding the object approximately every  $18^\circ$  plus

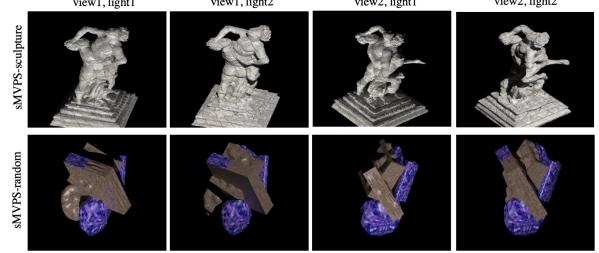


Figure 2. Example images from proposed synthetic MVPS dataset.

random jitter in position. For each view we render 10 randomly chosen directional light sources sampled uniformly on a  $45^\circ$  spherical cap centered at the camera's optical axis. In addition to the images, we render ground truth normals, depth, albedo, and roughness.

### 3.6. Training MVPSNet

We train MVPSNet with supervised loss over surface normal and depth using the ground-truth created with synthetic sMVPS dataset. For each reference view, we use 2 source views. And we randomly choose 3 lights out of 96 to train our model. The total loss is defined as:

$$L_{\text{mvps}} = \lambda_d \cdot Loss_d + \lambda_n \cdot Loss_n, \quad (10)$$

$$L_d = \sum_{s=1}^3 \lambda_{ds} \cdot L_{ds}, \quad \forall s = 1, 2, 3 \quad (11)$$

where  $L_{ds}$  and  $L_n$  refer to the depth loss for scale  $s$  and normal loss, respectively. For loss weights, we set  $L_n = 1$  and  $L_d = 10$ . The weights of each scale,  $\lambda_{ds} = 1$ , for  $s = 1, 2, 3$ .

## 4. Experiments

**Dataset.** We evaluate our method and conduct ablation study on DiLiGenT-MV [38] dataset, which is the only benchmark dataset for MVPS tasks and widely used by all previous approaches. It contains images of 5 objects with diverse materials captured from 20 views. For each view, the object is illuminated by one of 96 calibrated point light sources at one time, which gives us 96 images with varying lighting conditions.

**Evaluation Metrics.** We evaluate the quality of recovered meshes using L1 Chamfer distance from PyTorch3D [37, 53] and F-score [36] with L2 distance and 1mm threshold distance. The distances in both metrics are computed between the vertices of two meshes. The unit is mm.

**Evaluation details.** Ground truth meshes and meshes of PJ16 [51] and LZ20 [38] are included in the DiLiGenT-MV dataset [38]. We thank the authors of BKW22 [32] and BKC22 [30] for providing us with their reconstructed meshes. For PS-NeRF [65] we extract meshes for each object using their released code and unscale it to the scale of the ground truth as suggested in their code. The code or reconstructed meshes of BKW23 [31] was not released, so we only include their reported F-score on L2 distance in Table 3. To compare with PS-Transformer [27], we get normal

Category	Per-scene optimization					Generalizable			
	Manual Effort		Standalone			Single-view PS	MVS	MVPS	
Method	PJ16 [51]	LZ20 [38]	BKW22 [32]	BKC22 [30]	PS-NeRF [65]	PS-Transformer [27]	CasMVSNet [23]- RT	TransMVSNet [16]- RT	Ours
BEAR	2.54	0.73	1.01	1.01	<b>0.76</b>	3.17	1.47	1.48	0.80
BUDDHA	1.12	0.97	2.68	1.15	<b>0.86</b>	4.09	1.26	1.10	1.07
COW	1.14	0.39	1.09	<u>0.76</u>	<b>0.75</b>	3.04	1.27	1.05	0.77
POT2	3.21	0.67	1.54	1.40	<b>0.76</b>	3.05	1.46	1.05	0.82
READING	1.30	0.66	1.97	0.84	0.92	3.60	<u>0.75</u>	0.76	<b>0.66</b>
AVERAGE	1.86	0.69	1.66	1.03	<b>0.81</b>	3.39	1.24	1.09	0.82
Recon. Time/object	-	-	7 hrs	?	12 hrs	?	<b>22s</b>	<u>52s</u>	105s

Table 2. L1 Chamfer Distance (lower is better) between reconstructed mesh and GT after ICP. ‘-RT’ denotes trained on our synthetic MVPS dataset. For non-manual methods, the best result is shown in bold, 2nd best as underline. LZ20 & PJ16 involve carefully crafted steps, manual efforts in finding correspondence, and an initial mesh or point cloud.

Category	Per-scene optimization					Generalizable			
	Manual Effort		Standalone			Single-view PS	MVS	MVPS	
Method	PJ16 [51]	LZ20 [38]	BKW22 [32]	BKC22 [30]	BKW23* [31]	PS-Transformer [27]	CasMVSNet [23]-RT	TransMVSNet [16]-RT	Ours
BEAR	0.551	0.986	0.928	0.934	0.965	<b>0.995</b>	0.078	0.911	0.882
BUDDHA	0.940	0.936	0.687	0.926	<b>0.993</b>	<u>0.983</u>	0.066	0.919	0.963
COW	0.918	0.990	0.937	0.986	<u>0.987</u>	0.986	0.140	0.914	0.941
POT2	0.484	0.985	0.909	0.889	<u>0.991</u>	<u>0.991</u>	0.101	0.901	0.964
READING	0.905	0.975	0.810	0.971	0.975	0.961	0.961	<u>0.980</u>	0.978
AVERAGE	0.760	0.974	0.854	0.941	0.982	<u>0.983</u>	0.269	0.925	0.946
Recon. Time/object	-	-	7 hrs	?	?	12 hrs	?	<b>22s</b>	<u>52s</u>

Table 3. F-score with L2 distance and 1mm threshold distance (higher is better) between reconstructed mesh and GT after ICP. ‘-RT’ denotes trained on our synthetic MVPS dataset. For non-manual methods, the best result is shown in bold, 2nd best as underline. LZ20 & PJ16 involve carefully crafted steps, manual efforts in finding correspondence, and an initial mesh or point cloud. BKW23\* code not available, results from the paper.

maps from their pretrained model and integrate normals into depth maps, followed by a depth fusion step after rescaling all depth maps to the ground truth depth scales. To compare with CasMVSNet [23] and TransMVSNet [16], we consider both the pretrained models on DTU dataset [2] and models retrained on our synthetic dataset, dubbed as CasMVSNet-RT and TransMVSNet-RT. Images from 5 views are used to generate each single-view depth map and all 20 depth maps are fused together. For ours, we take images from 5 views and 10 lightings conditions for each view, along with corresponding light directions, as input to generate single-view depth maps and all 20 depth maps are fused for full recovery. Since there is no image showing the bottoms of objects in the dataset, following BKW22 [32] and BKC22 [30], we remove points that are located lower than +5 on the z-axis from all reconstructed meshes and the ground truth. Similar to most previous approaches [30, 32, 38, 51] we also perform a rigid registration using Iterative Closest Point (ICP) [3, 5, 13, 72] between the ground-truth and each reconstructed mesh for a fair comparison.

#### 4.1. Comparison with Existing Approaches.

We compare our algorithm with approaches that require per-scene optimization or training and with feed-forward generalizable methods. The quantitative result is showed in Table 2 and 3. We also show visual comparison of meshes from representative methods in Figure 3.

Method	CasMVSNet RT	CasMVSNet- RT	Ours (train 1 light/view)	Ours (train 3 light/view)
BEAR	2.00	1.47	1.31	0.80
BUDDHA	1.44	1.26	1.26	1.07
COW	2.73	1.27	1.06	0.77
POT2	1.89	1.46	1.07	0.82
READING	1.07	0.75	0.77	0.66
AVERAGE	1.83	1.24	1.09	0.82

Table 4. Ablation: Results are improved by retraining CasMVSNet [22] on proposed synthetic dataset (sMVPS). LAFM help in aggregating features across lighting variations, Ours (train 3 light/view) vs Ours (train 1 light/view), and is more accurate than CasMVSNet features, Ours (train 1 light/view) vs CasMVSNet-RT.

(i) **Per-scene Optimization.** Per-Scene optimization methods can also be categorized into:

(a) **Manual Efforts Needed.** We compare with two traditional multi-stage MVPS methods, PJ16 [51] and LZ20 [38], which require manual efforts. We outperform PJ16 [51] with a clear margin. Although LZ20 [38] achieves better results than ours, note that both methods, PJ16 & LZ20, consist of multiple steps with carefully crafted geometric modeling. Besides, they require an initial mesh [51] or point cloud [38] to build upon and their performance is sensitive to the initialization quality. When the initialization step fails in large textureless regions, LZ20 [38] incor-

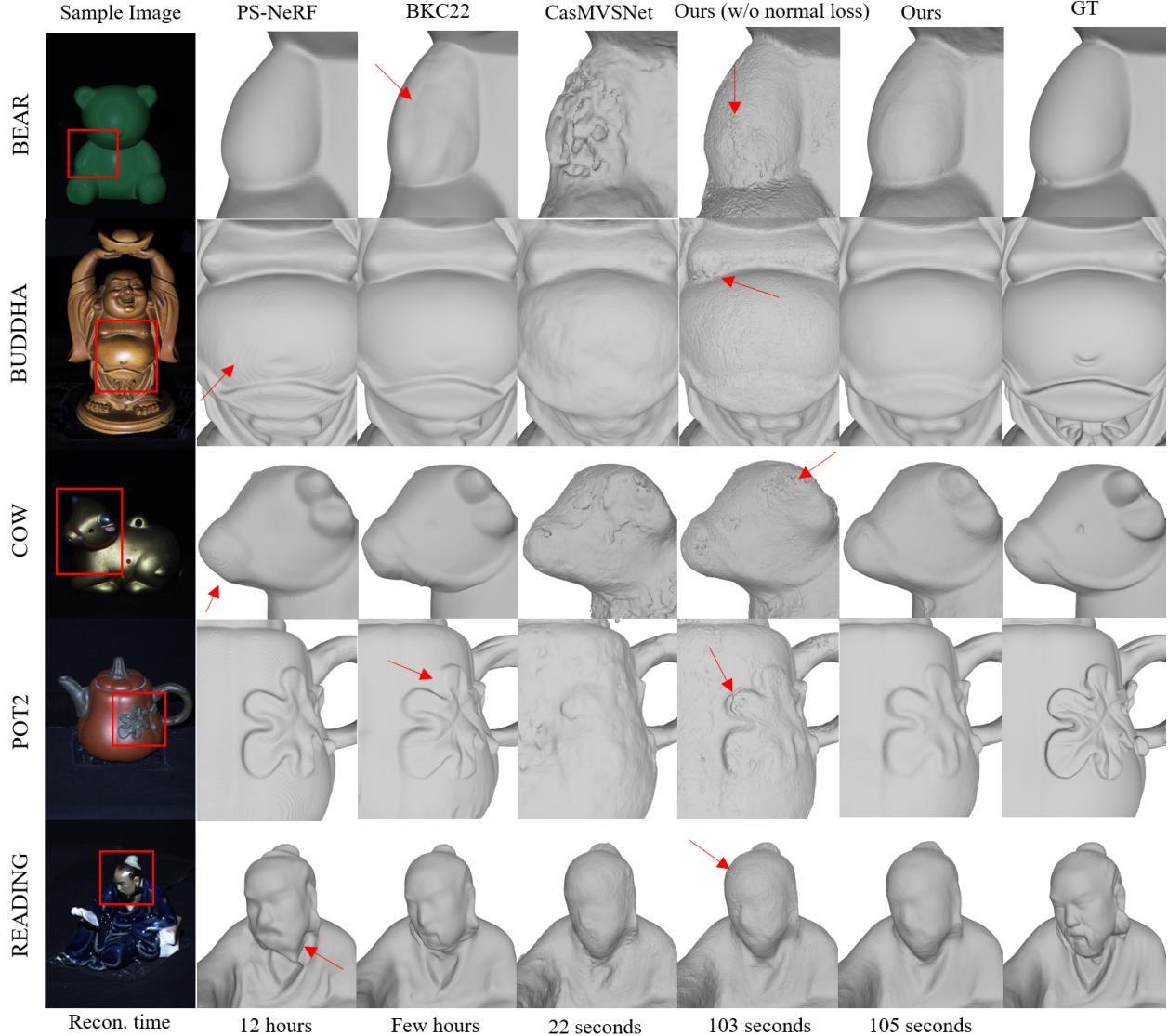


Figure 3. Qualitative comparison of our method with existing approaches (red arrow highlights artifacts) and ablation studies.

porates manual labeling to establish correspondence across views. In contrast, our pipeline is completely automatic without any manual efforts, does not involve any carefully crafted multi-stage approach, and does not require separate hyper-parameters for individual object.

**(b) Standalone Methods.** Recent deep learning-based MVPS methods, including BKW22 [30], BKC22 [32], BKW23 [31] and PS-NeRF [65], are simpler and easier to adopt. However, like traditional methods [38, 51], they still optimize one model for each object individually, resulting in low computational efficiency. In contrast, although our model is trained only on synthetic data, we outperform some per-scene optimized methods [30, 32] and get comparable results as state-of-the-art, PS-NeRF [65]. Furthermore, even though PS-NeRF [65] recovers high-quality meshes with details, its recovered surfaces contain

iso-contour pattern artifacts, *e.g.* see the red arrows in Fig. 3 on BUDDHA and COW, and often incorrect shapes, *e.g.* READING. Note that, we could not report L1 Chamfer distance on BKW23 [31], since the code is unavailable, but we show in Table 3 that our method is slightly better than BKW23 in F-score.

(ii) **Generalizable.** We compare our method with two categories of generalizable methods:

**(a) Single-view Photometric Stereo.** PS-Transformer [27] is a state-of-the-art PS network, which takes multiple images with the same viewpoint but different lighting conditions as input and generates single-view normal map prediction. To get 3D reconstruction, we integrate each normal map into a depth map and fuse them together. Since the integrated depths are of arbitrary scale, we rescale them to the range of ground truth depth. Inherently, PS methods

struggle with global shape modeling and it is challenging to stitch multi-view integrated depth of arbitrary scale, so PS-Transformer [27] doesn't perform well on full-view recovery.

**(b) Multi-view Stereo.** We also compare our method with CasMVSNet [23] and TransMVSNet [16], which takes a single lighting image for each view. For fairness, we re-train both methods using our synthetic dataset with suggested hyper-parameters in original papers . We observe that lighting information can largely improve both accuracy and quality, quantitatively and qualitatively. On textureless objects, *e.g.* BEAR, and textureless regions, *e.g.* belly of BUDDHA, MVS alone gets noisy and rough surfaces. Moreover, our meshes have more high-frequency details that MVS alone may struggle with, *e.g.* the texture on POT2. This is because our LAFMs are supervised with normal maps so they can learn the high-frequency components.

## 4.2. Computational efficiency

While our method outperforms some per-scene optimized methods [30, 32] and produces comparable results to state-of-the-arts [31, 65], the key advantage of our method is that it is fast, generalizable, and computationally efficient. Thus we analyze to the best of our abilities the inference time of various MVPS algorithms compared in this paper.

- LZ20 [38]: This algorithm takes *117 minutes* per object, without considering the time required for initializing a point cloud or any manual efforts.
- BKW22 [32]: takes 7 hours to train per object.
- BKC22 [30], BKW23 [31]: Since the authors did not mention the time required for training these algorithms it is not possible to provide an exact estimate. However, these approaches are based on MLPs, which take hours to train.
- PS-NeRF [65]: takes *12 hours* to train per object.
- CasMVSNet [22]:, in contrast, takes only 22 seconds per object, including obtaining depth maps for each view using 5 views (1 reference view and 4 source views) and 1 lighting per view (5 images processed for estimating a depth map), fusing depth maps from all 20 views to a point cloud, computing normals for vertices in the point cloud and adopting Screened Poisson [34] to recover a mesh from the point cloud.
- MVPSNet (Ours): takes a total of *105 seconds* to create a mesh, including steps of obtaining depth maps for each view using 5 views (1 reference view and 4 source views) and 10 lightings per view (50 images processed in total), fusing depth maps from all 20 views to a point cloud, and adopting Screened Poisson [34] to recover a mesh from the point cloud.

In summary, we are around  $240\times$  faster than BKW22 [32] and around  $411.4\times$  faster than PS-NeRF [65] ignoring their mesh extraction time.

## 4.3. Ablation study

The key contribution of this work includes: (a) our synthetic MVPS dataset, sMVPS-sculpture and sMVPS-random, and (b) ‘Light Augmented Feature Maps’ (LAFM)

for more accurate mesh reconstruction. Here we design experiments to analyze impacts of these contributions.

**Synthetic Data (sMVPS-sculpture and sMVPS-random).** To illustrate the effectiveness of our synthetic data, we evaluate the performance of a CasMVSNet [23] model trained on DTU dataset [2], and compare it with the CasMVSNet-RT model trained on our sMVPS dataset. The L1 Chamfer distance metric is reported in Table 4. We observe that training on our synthetic dataset improves reconstruction quality by 32.2%, proving the effectiveness of our proposed data for MVPS reconstructions. See supplementary materials for comparison between pretrained TransMVSNet [16] and TransMVSNet-RT trained on our synthetic dataset.

**Light Augmented Feature Maps (LAFM).** LAFM plays two key roles in our approach: (i) it aggregates features from images captured with multiple lighting conditions but the same viewpoints, and (ii) it is trained with surface normal loss which helps to preserve high-frequency details in the features. The predicted surface normals are used to further refine the reconstructed mesh.

For understanding the impact of (i), we train our proposed MVPSNet with just a single lighting image per view instead of 3 lightings. Thus LAFMs are only aggregated across 1 image per view. In Table 4 we observe that using a single lighting per view (‘Ours (train 1 light/view)’) produces worse results (1.09 vs 0.82) than using 3 lightings per view (‘Ours (train 3 light/view)’). However, even a single lighting per view produces better performance than CasMVSNet-RT (1.09 vs 1.24), which is also trained on single lighting per view. This shows that LAFM is effective in both extracting accurate information from just a single image and aggregating shading information across multiple images with varying illumination.

For understanding the impact of (ii), we train our proposed MVPSNet without any surface normal loss, ‘Ours (w/o normal loss)’. We observe ‘Ours (w/o normal loss)’ is quantitatively comparable to ‘Ours (w/normal loss)’, 0.79 vs 0.82 in L1 chamfer distance and 0.985 vs 0.985 in F-score. However, in Fig. 3 we observe that the meshes produced by ‘Ours (w/o normal loss)’ are significantly noisier as shown with red arrows.

## 5. Conclusion

In this work, we propose a fast and generalizable approach for MVPS. We introduce Light Aggregated Feature Maps that leverage shading cues from images with the same view under multiple lighting conditions to produce richer features in textureless regions. Being trained with normal estimation, LAFM also enable higher quality reconstruction than traditional MVS methods with only little compromise to speed. When trained on the synthetic sMVPS dataset we propose, our method produces results comparable to STOA method that is about 400x slower at inference time.

## References

- [1] 3d textures. <https://3dtextures.me/>. Accessed: 2020. 12
- [2] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120:153–168, 2016. 6, 8, 16
- [3] K. S. Arun, T. S. Huang, and S. D. Blostein. Least-squares fitting of two 3-d point sets. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PAMI-9(5):698–700, 1987. 6, 12, 14
- [4] Meghna Asthana, William AP Smith, and Patrik Huber. Neural apparent brdf fields for multiview photometric stereo. *arXiv preprint arXiv:2207.06793*, 2022. 3
- [5] Paul J Besl and Neil D McKay. Method for registration of 3-d shapes. In *Sensor fusion IV: control paradigms and data structures*, volume 1611, pages 586–606. Spie, 1992. 6, 12, 14
- [6] Michael Bleyer, Christoph Rhemann, and Carsten Rother. Patchmatch stereo-stereo matching with slanted support windows. In *Bmvc*, volume 11, pages 1–11, 2011. 2
- [7] Mark Boss, Varun Jampani, Kihwan Kim, Hendrik P.A. Lensch, and Jan Kautz. Two-shot spatially-varying brdf and shape estimation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [8] Chenjie Cao, Xinlin Ren, and Yanwei Fu. Mvsformer: Learning robust image representations via transformers and temperature-based depth for multi-view stereo. *arXiv preprint arXiv:2208.02541*, 2022. 1
- [9] Guanying Chen, Kai Han, Boxin Shi, Yasuyuki Matsushita, and Kwan-Yee K Wong. Self-calibrating deep photometric stereo networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8739–8747, 2019. 2, 4
- [10] Guanying Chen, Kai Han, Boxin Shi, Yasuyuki Matsushita, and Kwan-Yee Kenneth Wong. Deep photometric stereo for non-lambertian surfaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [11] Guanying Chen, Kai Han, and Kwan-Yee K Wong. Ps-fcn: A flexible learning framework for photometric stereo. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–18, 2018. 2
- [12] Guanying Chen, Michael Waechter, Boxin Shi, Kwan-Yee K Wong, and Yasuyuki Matsushita. What is learned in deep uncalibrated photometric stereo? In *European Conference on Computer Vision*, 2020. 2
- [13] Yang Chen and Gérard Medioni. Object modelling by registration of multiple range images. *Image and vision computing*, 10(3):145–155, 1992. 6, 12, 14
- [14] Robert T Collins. A space-sweep approach to true multi-image matching. In *Proceedings CVPR IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 358–363. Ieee, 1996. 2
- [15] Yikang Ding, Wentao Yuan, Qingtian Zhu, Haotian Zhang, Xiangyue Liu, Yuanjiang Wang, and Xiao Liu. Transmvsnet: Global context-aware multi-view stereo network with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8585–8594, 2022. 1
- [16] Yikang Ding, Wentao Yuan, Qingtian Zhu, Haotian Zhang, Xiangyue Liu, Yuanjiang Wang, and Xiao Liu. Transmvsnet: Global context-aware multi-view stereo network with transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8585–8594, 2022. 2, 6, 8, 12, 13, 14, 16
- [17] Yasutaka Furukawa, Brian Curless, Steven M Seitz, and Richard Szeliski. Towards internet-scale multi-view stereo. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pages 1434–1441. IEEE, 2010. 2
- [18] Yasutaka Furukawa and Jean Ponce. Accurate, dense, and robust multiview stereopsis. *IEEE transactions on pattern analysis and machine intelligence*, 32(8):1362–1376, 2009. 2
- [19] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Massively parallel multiview stereopsis by surface normal diffusion. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 873–881, 2015. 14
- [20] Khang Truong Giang, Soohwan Song, and Sungho Jo. Curvature-guided dynamic scale networks for multi-view stereo. *arXiv preprint arXiv:2112.05999*, 2021. 1
- [21] Michael Goesele, Noah Snavely, Brian Curless, Hugues Hoppe, and Steven M Seitz. Multi-view stereo for community photo collections. In *2007 IEEE 11th International Conference on Computer Vision*, pages 1–8. IEEE, 2007. 2
- [22] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. 2019. 1, 6, 8, 12, 13, 14, 16
- [23] Xiaodong Gu, Zhiwen Fan, Siyu Zhu, Zuozhuo Dai, Feitong Tan, and Ping Tan. Cascade cost volume for high-resolution multi-view stereo and stereo matching. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2495–2504, 2020. 2, 3, 4, 6, 8, 15
- [24] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 12
- [25] Carlos Hernandez, George Vogiatzis, and Roberto Cipolla. Multiview photometric stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(3):548–554, 2008. 2
- [26] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2821–2830, 2018. 2
- [27] Satoshi Ikehata. Ps-transformer: Learning sparse photometric stereo network using self-attention mechanism. *arXiv preprint arXiv:2211.11386*, 2022. 2, 5, 6, 7, 8, 15
- [28] Mengqi Ji, Juergen Gall, Haitian Zheng, Yebin Liu, and Lu Fang. Surfacenet: An end-to-end 3d neural network for multiview stereopsis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2307–2315, 2017. 2

- [29] Sing Bing Kang, Richard Szeliski, and Jinxiang Chai. Handling occlusions in dense multi-view stereo. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 1, pages I–I. IEEE, 2001. 2
- [30] Berk Kaya, Suryansh Kumar, Carlos Oliveira, Vittorio Ferrari, and Luc Van Gool. Uncertainty-aware deep multi-view photometric stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12601–12611, 2022. 1, 2, 3, 5, 6, 7, 8, 15
- [31] Berk Kaya, Suryansh Kumar, Carlos Oliveira, Vittorio Ferrari, and Luc Van Gool. Multi-view photometric stereo revisited. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 3126–3135, 2023. 1, 2, 3, 5, 6, 7, 8, 15
- [32] Berk Kaya, Suryansh Kumar, Francesco Sarno, Vittorio Ferrari, and Luc Van Gool. Neural radiance fields approach to deep multi-view photometric stereo. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1965–1977, 2022. 1, 2, 5, 6, 7, 8, 14, 15
- [33] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, volume 7, page 0, 2006. 2
- [34] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 2, 5, 8, 14
- [35] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 13
- [36] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 5, 14
- [37] Christoph Lassner and Michael Zollhöfer. Pulsar: Efficient sphere-based neural rendering. *arXiv:2004.07484*, 2020. 5
- [38] Min Li, Zhenglong Zhou, Zhe Wu, Boxin Shi, Changyu Diao, and Ping Tan. Multi-view photometric stereo: A robust solution and benchmark dataset for spatially varying isotropic materials. *IEEE Transactions on Image Processing*, 29:4159–4173, 2020. 1, 2, 5, 6, 7, 8, 13, 14, 15, 16, 17, 18, 19
- [39] Zhengqin Li, Zexiang Xu, Ravi Ramamoorthi, Kalyan Sunkavalli, and Manmohan Chandraker. Learning to reconstruct shape and spatially-varying reflectance from a single image. In *SIGGRAPH Asia 2018 Technical Papers*, page 269. ACM, 2018. 2
- [40] Daniel Lichy, Soumyadip Sengupta, and David W Jacobs. Fast light-weight near-field photometric stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12612–12621, 2022. 2
- [41] Daniel Lichy, Soumyadip Sengupta, and David W. Jacobs. Fast light-weight near-field photometric stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12612–12621, June 2022. 2
- [42] Daniel Lichy, Jiaye Wu, Soumyadip Sengupta, and David W. Jacobs. Shape and material capture at home. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6123–6133, June 2021. 2
- [43] Daniel Lichy, Jiaye Wu, Soumyadip Sengupta, and David W Jacobs. Shape and material capture at home. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6123–6133, 2021. 5, 12
- [44] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987. 2
- [45] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. P-mvsnet: Learning patch-wise matching confidence aggregation for multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10452–10461, 2019. 2
- [46] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 2
- [47] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 2
- [48] Diego Nehab, Szymon Rusinkiewicz, James Davis, and Ravi Ramamoorthi. Efficiently combining positions and normals for precise 3d geometry. *ACM transactions on graphics (TOG)*, 24(3):536–543, 2005. 2
- [49] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5589–5599, 2021. 3
- [50] Jaesik Park, Sudipta N Sinha, Yasuyuki Matsushita, Yu-Wing Tai, and In So Kweon. Multiview photometric stereo using planar mesh parameterization. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1161–1168, 2013. 2
- [51] Jaesik Park, Sudipta N Sinha, Yasuyuki Matsushita, Yu-Wing Tai, and In So Kweon. Robust multiview photometric stereo using planar mesh parameterization. *IEEE transactions on pattern analysis and machine intelligence*, 39(8):1591–1604, 2016. 1, 2, 5, 6, 7, 15
- [52] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 12
- [53] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 5
- [54] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Nether-*

- lands, October 11-14, 2016, Proceedings, Part III* 14, pages 501–518. Springer, 2016. 2
- [55] Soumyadip Sengupta, Angjoo Kanazawa, Carlos D Castillo, and David W Jacobs. Sfsnet: Learning shape, reflectance and illuminance of faces in the wild’. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6296–6305, 2018. 2
- [56] Boxin Shi, Zhe Wu Mo, Dinglong Duan, Sai-Kit Yeung, and Ping Tan. A benchmark dataset and evalution for non-lambertian and uncalibrated photometric stereo. *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 41(2):271–284, 2019. 2
- [57] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7495–7504, 2021. 1
- [58] Christoph Strecha, Rik Fransens, and Luc Van Gool. Wide-baseline stereo from multiple views: a probabilistic account. In *Proceedings of the 2004 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2004. CVPR 2004.*, volume 1, pages I–I. IEEE, 2004. 2
- [59] Christoph Strecha, Rik Fransens, and Luc Van Gool. Combined depth and outlier estimation in multi-view stereo. In *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, volume 2, pages 2394–2401. IEEE, 2006. 2
- [60] Jiaming Sun, Yiming Xie, Linghao Chen, Xiaowei Zhou, and Hujun Bao. NeuralRecon: Real-time coherent 3D reconstruction from monocular video. *CVPR*, 2021. 1
- [61] Olivia Wiles and Andrew Zisserman. Silnet : Single- and multi-view reconstruction by learning from silhouettes. In *British Machine Vision Conference 2017, BMVC 2017, London, UK, September 4-7, 2017*. BMVA Press, 2017. 2, 5, 12
- [62] Robert J Woodham. Photometric method for determining surface orientation from multiple images. *Optical engineering*, 19(1):139–144, 1980. 2
- [63] Qingshan Xu and Wenbing Tao. Learning inverse depth regression for multi-view stereo with correlation cost volume. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 12508–12515, 2020. 2
- [64] Zexiang Xu, Kalyan Sunkavalli, Sunil Hadap, and Ravi Ramamoorthi. Deep image-based relighting from optimal sparse samples. *ACM Transactions on Graphics (TOG)*, 37(4):126, 2018. 2, 5, 12
- [65] Wenqi Yang, Guanying Chen, Chaofeng Chen, Zhenfang Chen, and Kwan-Yee K Wong. Ps-nerf: Neural inverse rendering for multi-view photometric stereo. *arXiv preprint arXiv:2207.11406*, 2022. 1, 2, 3, 5, 6, 7, 8, 15, 20
- [66] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Proceedings of the European conference on computer vision (ECCV)*, pages 767–783, 2018. 1, 2
- [67] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. *Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [68] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Basri Ronen, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. *Advances in Neural Information Processing Systems*, 33:2492–2502, 2020. 1
- [69] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images, 2020. 2
- [70] Kai Zhang, Fujun Luan, Qianqian Wang, Kavita Bala, and Noah Snavely. Physg: Inverse rendering with spherical gaussians for physics-based material editing and relighting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5453–5462, 2021. 3
- [71] Xiuming Zhang, Pratul P Srinivasan, Boyang Deng, Paul Debevec, William T Freeman, and Jonathan T Barron. Nefactor: Neural factorization of shape and reflectance under an unknown illumination. *ACM Transactions on Graphics (TOG)*, 40(6):1–18, 2021. 1
- [72] Zhengyou Zhang. Iterative point matching for registration of free-form curves and surfaces. *International journal of computer vision*, 13(2):119–152, 1994. 6, 12, 14
- [73] Zhenglong Zhou, Zhe Wu, and Ping Tan. Multi-view photometric stereo with spatially varying isotropic materials. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1482–1489, 2013. 2

## 6. Appendix

### 6.1. Overview

In this supplementary material, we will include the following contents:

- We describe more details about our **sMVPS dataset** in **Section 2** and show additional example images in Figure 4 and Figure 5.
- We provide additional **experiment details**, including notations we use for network architecture and implementation details in **Section 3**.
- We explain our **mesh extraction pipeline** in detail in **Section 4** together with the parameters we use.
- We provide the equations of the **evaluation metrics** we use in **Section 5**.
- In the main paper, we provide L1 Chamfer distance and F-score with L2 distance after ICP [3, 5, 13, 72]. Here in **Section 6**, we also provide **results of L1 Chamfer distance and F-score with L2 distance before ICP** in Table 5 and 6.
- To further compare with generalizable MVS methods, in addition to CasMVSNet [22] in the main paper, we also **compare our method with a state-of-the-art MVS method, TransMVSNet [16]**, in **Section 7**. TransMVSNet [16] is built upon CasMVSNet and adopts a transformer to consider intra-image and inter-image feature interactions, which further improves the result of CasMVSNet [22]. See Table 7 for the comparison.
- We include additional qualitative results. We show the global shape of reconstructed mesh from each method under three different views in Figure 6 - 10. We also show additional zoomed areas for visual comparison between meshes in Figure 11.

### 6.2. sMVPS datasets

**Object and Camera Positioning** For both synthetic datasets objects are placed at the center of world coordinates with the object's up direction along the z-axis. Objects are scaled to be inside a sphere of radius one. We use a pinhole camera for rendering with an FOV of  $9.3^\circ$ , which is similar to the FOV used to capture the DiligentMV dataset. Camera positions are most easily described in spherical coordinates i.e. an azimuth angle, a polar angle, and a radial distance. The azimuth angle for the  $i$ th camera is  $(18 + X_i)^\circ$  Where  $X_i$  is a uniform random number between -3 and 3, and  $i$  runs from 0 to 19. The polar angle for each camera is sampled uniformly from  $62\text{-}64^\circ$ . The radial distance is sampled

uniformly between 14 and 16.5. This distance is chosen so the object occupies the majority of the image.

**Light Positioning** Each view is rendered under 10 directional lights. The first light is always co-directional with the camera while the other 9 are randomly sampled from the spherical cap centered on the cameras optical axis with an angle of  $45^\circ$ .

**BRDF** To generate BRDFs we follow [43]. Namely we use the Cook-Torrance BRDF model with spatially-varying albedo drawn from 415 free textures from [1], and randomly generated roughness as described in [43]. Roughness is constant in the case of sMVPS-sculpture and constant for each primitive in the case of sMVPS-random.

**Object Meshes** For the sMVPS-random dataset objects are drawn from the collections of random primitives generated by [64] using a 90-10 train/test split. For the sMVPS-sculpture dataset we use the following meshes from [61] to render the train set: nymph-seated, standing-isis-priest, the-slave-girl, thor, three-danish-polar-explorers, tiger-devouring-a-gavial, two-wrestlers-in-combat, ugolino-and-his-sons, virgin-and-child, woman-associated-with-the-cult-of-isis, wounded-amazon, wounded-cupid, wrestling-decimated-cleaned and the mesh virgin-mary-with-her-dead-son for the test set.

**Rendering** Images are rendered with Mitsuba 2 using the path-tracer integration method. We render at a resolution of 612x512 with 128 samples-per-pixel.

**More Examples** To further show the diversity on surface shape, texture and material of our sMVPS datasets, we provide additional example images of sMVPS-sculpture in Figure 4 and sMVPS-random in Figure 5.

### 6.3. Additional experiment details

#### 6.3.1 Notation in Figure 1 of main paper

The architecture of our network is illustrated in Figure 1 in the paper and we describe a few details and notations we use:

**ResBlk:** Resnet block. It consists of  $\text{conv2d}(\text{kernel}=3) \rightarrow \text{BatchNorm} \rightarrow \text{ReLU} \rightarrow \text{conv2d}(\text{kernel}=3) \rightarrow \text{BatchNorm}$ . And the input of this block is added to the output of this block as a residual connection [24].

**Tconv:** ConvTranspose2d layer in Pytorch with kernel=3.

#### 6.3.2 Implementation details

Our model is implemented in Pytorch [52] and we use a NVIDIA RTX A6000 GPU to train it. For input images, we crop them to  $512 \times 512$  and rescale the pixel values to  $(0, 1)$ . For each training sample, we use 3 views and 3 lightings. It is challenging to find correspondences for view selection in textureless regions, so we simply take the two adjacent views of a reference view as source views. To

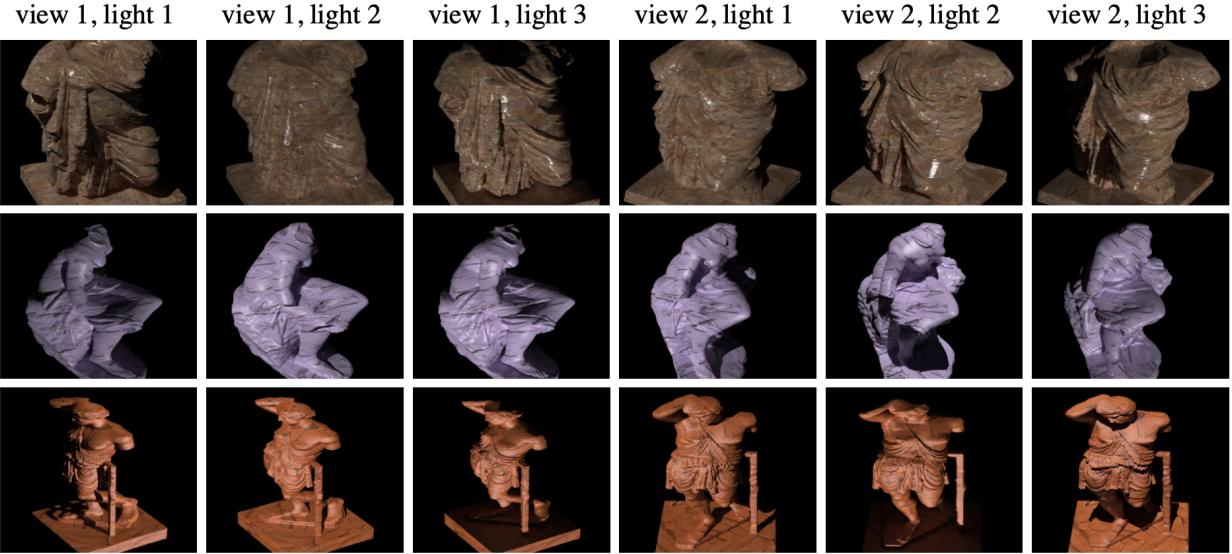


Figure 4. Additional example images of sMVPS-sculpture.

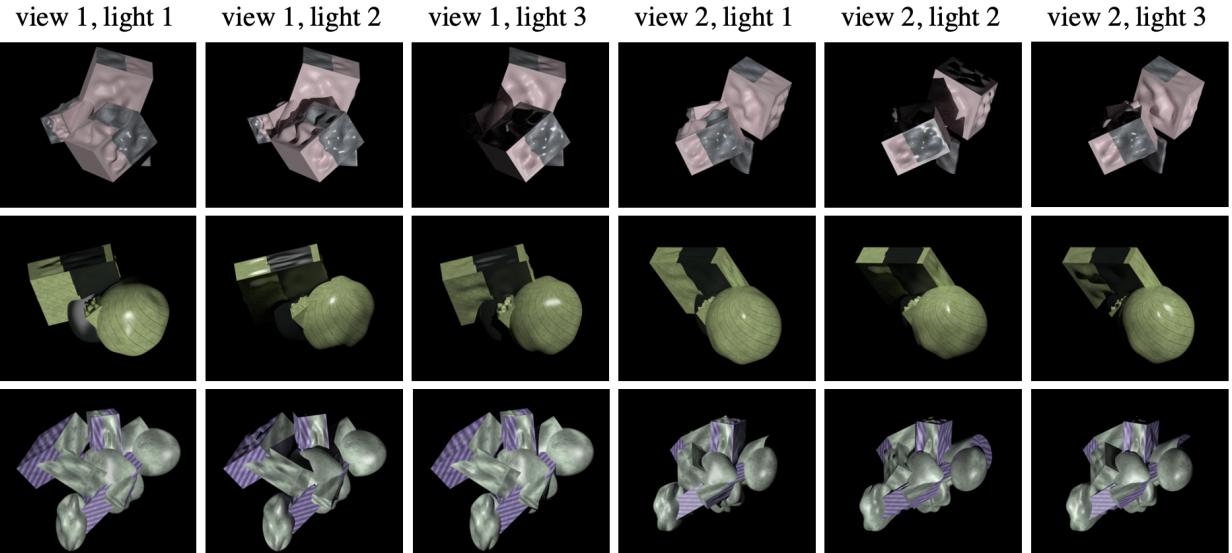


Figure 5. Additional example images of sMVPS-random.

make our model more robust to different lighting configurations, we randomly sample 3 lightings and use the same lightings for all views, resulting in  $3 \times 3 = 9$  images for each training sample. We use Adam [35] optimizer and set betas as  $(0.9, 0.999)$ . We trained 50 epochs in total. The initial learning rate is 0.001 and it decays to half at steps [8, 12, 30, 40]. To get ground truth depth map of DiLiGent-MV [38], we render depth map from ground truth mesh and camera parameters.

## 6.4. Mesh extraction pipeline

We use the same mesh extraction pipeline to recover 3D mesh from predicted depth maps for CasMVSNet [22], Ours and TransMVSNet [16] for a fair comparison.

### 6.4.1 Depth filtering

We use two kinds of masks to filter predicted depth maps. First, we employ 2D object mask to rule out background. This is because our model is only trained on pixels within

an object. Second, we apply geometric filtering to only keep depth predictions that are consistent across adjacent views. For each object pixel in reference view,  $p_o$ , we have a predicted depth aligned with this view,  $d_o$ . We lift  $p_o$  to a 3D point,  $P_o$ , and project  $P_o$  to a source view pixel,  $p_s$ . Assume the predicted depth of source view at  $p_s$  is  $d_s$ . By lifting  $p_s$  using  $d_s$ , we get a 3D point  $P_s$ . Projecting  $P_s$  back to reference view results in a reprojected pixel,  $p_r$  and a depth  $d_r$ . We set thresholds for the distance between the original pixel,  $p_o$ , and reprojected pixel,  $p_r$ , as well as relative difference between  $d_o$  and  $d_r$  as follows:

$$dist(p_o, p_r) < 1, \quad (12)$$

$$abs(d_o - d_r)/d_o < 0.01 \quad (13)$$

For  $p_o$  and  $d_o$ , we check its geometric consistency with each source view and keep it only if it is consistent with at least one source view.

#### 6.4.2 Depth fusion

After the depth filtering step, we combine each depth map in a fusion step. For an object pixel  $p_o$ , we simply average over  $d_o$  and all the estimations from source views that are consistent with it,  $d_{si}$  for  $i = 1, \dots, i_N$ , where  $i_N$  is the total number of geometric consistent neighboring views, and use this average as depth at  $p_o$ . We then lift  $p_o$  to a vertex in point cloud and attach the predicted normal,  $n_o$ , to it. This way, we get point cloud utilizing information from depth maps of all views.

Note there are other possible depth fusion methods, *e.g.* GIPUMA [19], some of which may achieve better fusion performance for certain datasets. But there is no method that is better for all datasets, so we leave exploration in this direction as a future work.

#### 6.4.3 Surface reconstruction

We apply Screend Poisson Surface Reconstruction (SPSR) [34] to recover mesh from point cloud. We use same set of parameters for all methods and all objects. Specifically, we set *reconstruction\_depth* = 8, *minimum\_number\_of\_samples* = 1.5 and *interpolation\_weight* = 4. Note that before recovering surfaces, an extra step of computing normal based on the point cloud is needed for methods without normal prediction, *i.e.*, CasMVSNet [22] and TransMVSNet [16].

#### 6.4.4 Evaluation metrics details

We use L1 Chamfer distance and F-score with L2 distance (threshold at 1mm) to evaluate the quality of reconstructed mesh. Both metrics are applied to two sets of 3D points, which are vertices of reconstructed mesh and ground truth mesh.

Give two point sets,  $\mathcal{R}$  and  $\mathcal{G}$ , L1 Chamfer distance is defined as follows:

$$CD(\mathcal{R}, \mathcal{G}) = \frac{1}{|\mathcal{R}|} \sum_{x \in \mathcal{R}} \min_{y \in \mathcal{G}} \|x - y\| + \frac{1}{|\mathcal{G}|} \sum_{y \in \mathcal{G}} \min_{x \in \mathcal{R}} \|x - y\|. \quad (14)$$

We use the F-score similarly defined as [36], for a reconstructed point  $r \in \mathcal{R}$ , its L2 distance to the ground truth mesh  $\mathcal{G}$  is

$$e_{r \rightarrow \mathcal{G}} = \min_{g \in \mathcal{G}} \|r - g\|_2, \quad (15)$$

and for a ground truth point  $g \in \mathcal{G}$ , its distance to the reconstructed mesh is defined as:

$$e_{g \rightarrow \mathcal{R}} = \min_{r \in \mathcal{R}} \|r - g\|_2, \quad (16)$$

The precision and recall for a threshold  $d$  are:

$$P(d) = \frac{1}{|\mathcal{R}|} \sum_{r \in \mathcal{R}} [e_{r \rightarrow \mathcal{G}} < d] \quad (17)$$

$$R(d) = \frac{1}{|\mathcal{G}|} \sum_{g \in \mathcal{G}} [e_{g \rightarrow \mathcal{R}} < d] \quad (18)$$

F-score is the harmonic mean of precision and recall as a summary measure:

$$F(d) = \frac{2P(d)R(d)}{P(d) + R(d)} \quad (19)$$

#### 6.5 Additional results without ICPs

In the paper, we report results after ICP [3, 5, 13, 72], which is an extra registration step we applied to all meshes after being reconstructed. It is initially aimed to fairly compare our mesh with others as several methods indicate that they did registration after extracting meshes [32, 38]. However, we find it also helpful to improve accuracy of other methods, even for those that already have registration applied. Since there is no standard way to do registration among existing methods, we applied ICP to meshes from all methods, regardless of whether they have done registration or not.

For a complete comparison, we also provide the quantitative results of L1 Chamfer distance and F-score with L2 distance (threshold at 1mm) without ICP [3, 5, 13, 72] in Table 5 and Table 6, respectively. They show that even without registration, our method can still perform comparably with state-of-the-art methods with registration.

#### 6.6 Comparison with TransMVSNet

TransMVSNet [16] is a state-of-the-art MVS method leveraging a transformer to extract both intra-image global context and inter-image feature interaction. We retrain a TransMVSNet model on our synthetic MVPS dataset using hyper-parameters suggested by [16]. Each training sample has 3 views and test sample has 5 views, which is the same set-up as our method and the retrained CasMVSNet [22].

We provide L1 Chamfer distance and F-score with L2 distance (threshold at 1mm) in Table 7, which demonstrate

Category	Per-scene optimization					Generalizable		
	Manual Effort		Standalone			Single-view PS	MVS	MVPS
Method	PJ16 [51]	LZ20 [38]	BKW22 [32]	BKC22 [30]	PS-NeRF [65]	PS-Transformer [27]	CasMVSNet [23]- RT	Ours
BEAR	2.63	0.74	1.03	1.09	<b>0.81</b>	3.25	1.38	<u>0.91</u>
BUDDHA	1.18	0.99	2.44	1.19	<b>0.98</b>	4.44	1.30	<u>1.12</u>
COW	1.16	0.39	1.08	0.86	<b>0.78</b>	2.67	1.26	<u>0.80</u>
POT2	3.27	0.69	1.32	1.32	<b>0.81</b>	2.92	1.43	<u>0.94</u>
READING	1.49	0.74	1.94	0.93	0.98	3.69	<u>0.83</u>	<b>0.76</b>
AVERAGE	1.95	0.71	1.56	1.08	<b>0.87</b>	3.39	1.24	<u>0.91</u>

Table 5. L1 Chamfer Distance (lower is better) between reconstructed mesh and GT without ICP. ‘-RT’ denotes trained on our synthetic MVPS dataset. For non-manual methods, the best result is shown in bold, 2nd best as underline. LZ20 & PJ16 involve carefully crafted steps, manual efforts in finding correspondence, and an initial mesh or point cloud.

Category	Per-scene optimization						Generalizable		
	Manual Effort		Standalone				Single-view PS	MVS	MVPS
Method	PJ16 [51]	LZ20 [38]	BKW22 [32]	BKC22 [30]	BKW23* [31]	PS-NeRF [65]	PS-Transformer [27]	CasMVSNet [23]-RT	Ours
BEAR	0.504	0.987	0.926	0.895	0.965	<b>0.994</b>	0.496	0.902	<u>0.990</u>
BUDDHA	0.935	0.935	0.745	0.922	<b>0.993</b>	<u>0.970</u>	0.387	0.913	0.953
COW	0.917	0.990	0.943	0.981	<u>0.987</u>	0.984	0.617	0.896	<b>0.993</b>
POT2	0.459	0.985	0.929	0.909	<u>0.991</u>	0.990	0.609	0.891	<b>0.992</b>
READING	0.868	0.975	0.807	0.970	<u>0.975</u>	0.946	0.501	<u>0.981</u>	<b>0.989</b>
AVERAGE	0.737	0.974	0.870	0.935	<u>0.982</u>	0.977	0.522	0.917	<b>0.983</b>

Table 6. F-score on L2 distance (higher is better) between reconstructed mesh and GT without ICP. ‘-RT’ denotes trained on our synthetic MVPS dataset. For non-manual methods, the best result is shown in bold, 2nd best as underline. LZ20 & PJ16 involve carefully crafted steps, manual efforts in finding correspondence, and an initial mesh or point cloud. BKW23\* code not available, result from the paper.

that our method outperform STOA generalizable MVS method. Also note that the transformer used in TransMVS-Net increases the runtime of the method compared to CasMVSNet.

Metrics	L1 Chamfer distance					F-score (1mm)				
	CasMVSNet [22]	CasMVSNet- RT	TransMVSNet [16]	TransMVSNet- RT	Ours	CasMVSNet [22]	CasMVSNet- RT	TransMVSNet [16]	TransMVSNet- RT	Ours
BEAR	2.00	1.47	1.02	1.48	<b>0.80</b>	0.789	0.911	0.962	0.882	<b>0.991</b>
BUDDHA	1.44	1.26	1.09	1.10	<b>1.07</b>	0.878	0.919	0.961	0.963	<b>0.958</b>
COW	2.73	1.27	1.15	1.05	<b>0.77</b>	0.658	0.914	0.927	0.941	<b>0.993</b>
POT2	1.89	1.46	1.10	1.05	<b>0.82</b>	0.799	0.901	0.956	0.964	<b>0.994</b>
READING	1.07	0.75	0.87	0.76	<b>0.66</b>	0.941	0.980	0.971	0.978	<b>0.988</b>
AVERAGE	1.83	1.24	1.05	1.09	<b>0.82</b>	0.813	0.925	0.955	0.946	<b>0.985</b>
Recon. Time/object	22s	22s	52s	52s	105s	22s	22s	52s	52s	105s

Table 7. Comparison with TransMVSNet on L1 Chamfer distance and F-score with L2 distance (threshold at 1mm) after ICP. CasMVSNet [22] and TransMVSNet [16] denote the pretrained models on DTU dataset [2]. 'RT' denotes trained on our synthetic MVPS dataset.

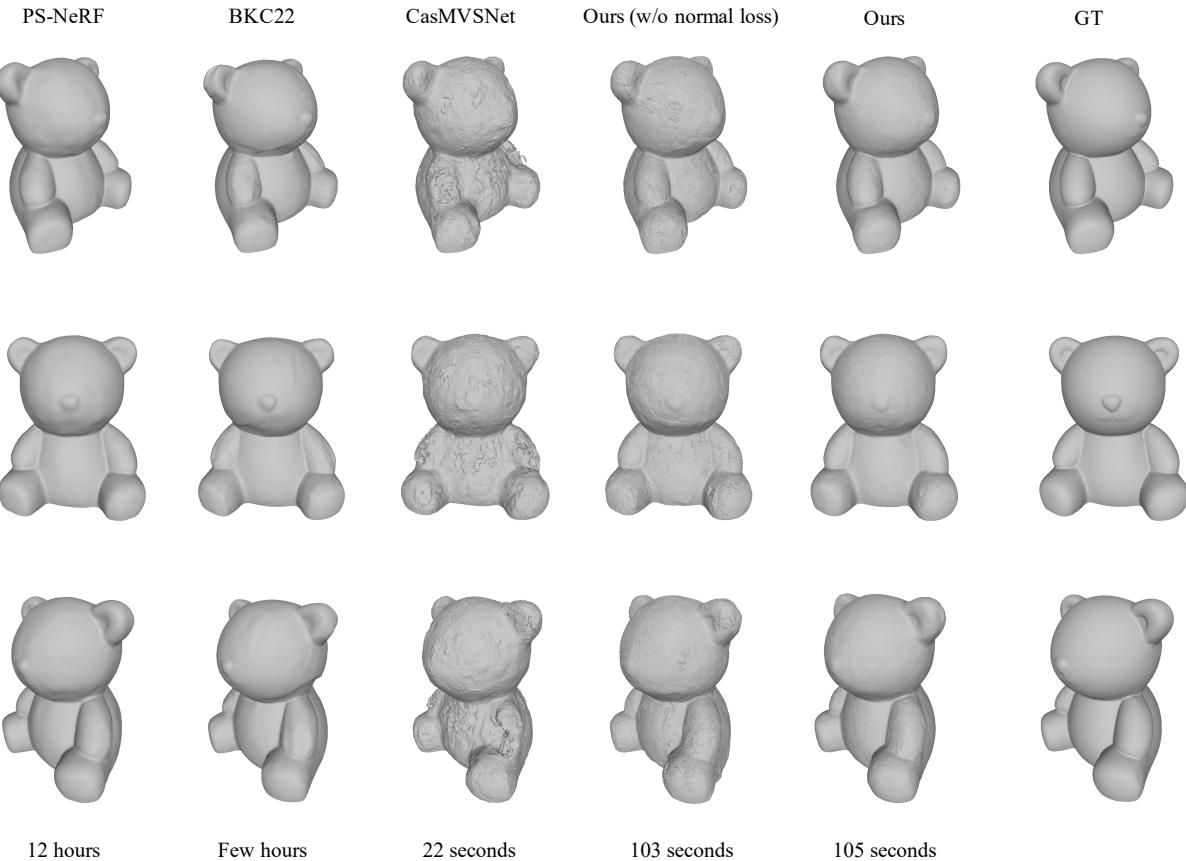


Figure 6. Reconstruction of BEAR under three different views (left-side, front, right-side) in DiLiGenT-MV [38]. Last row is reconstruction time.

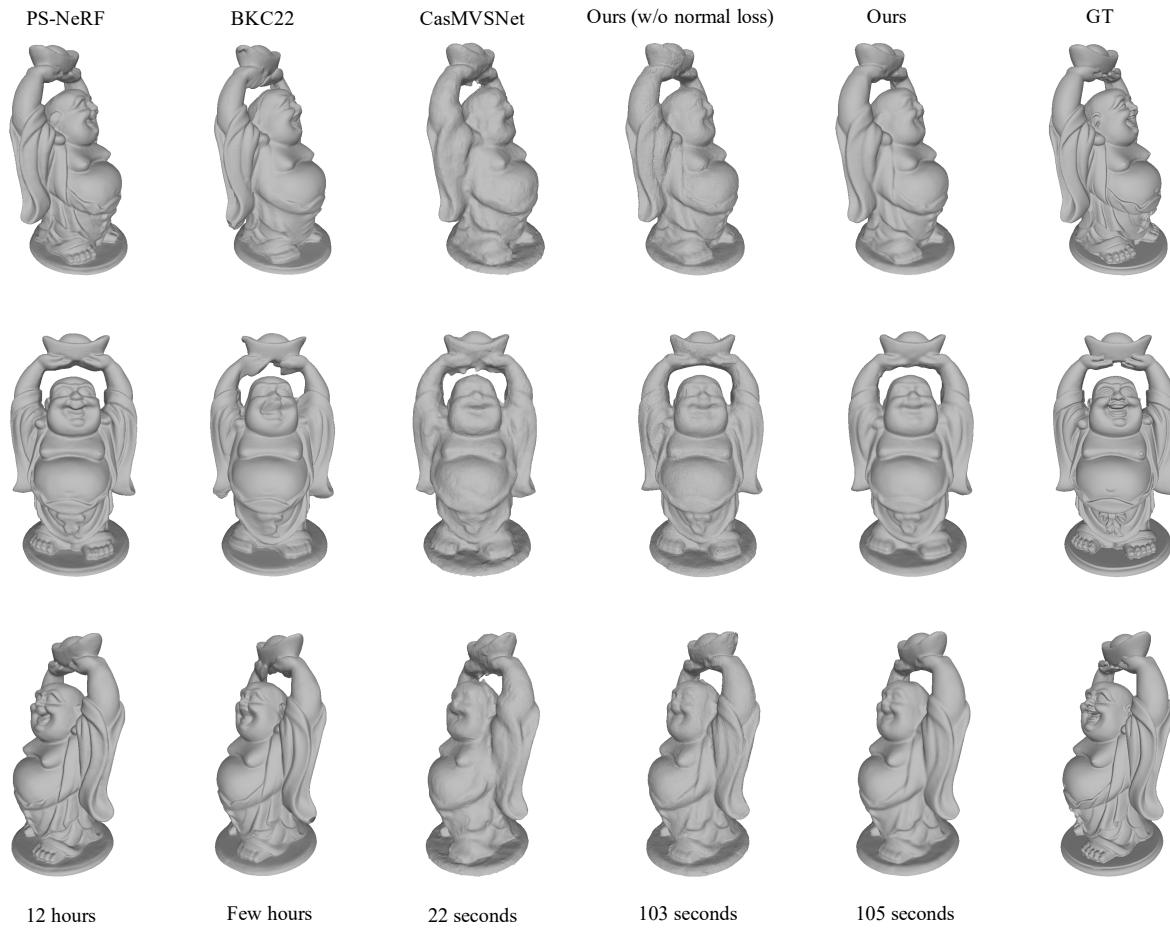


Figure 7. Reconstruction of BUDDHA under three different views (left-side, front, right-side) in DiLiGenT-MV [38]. Last row is reconstruction time.

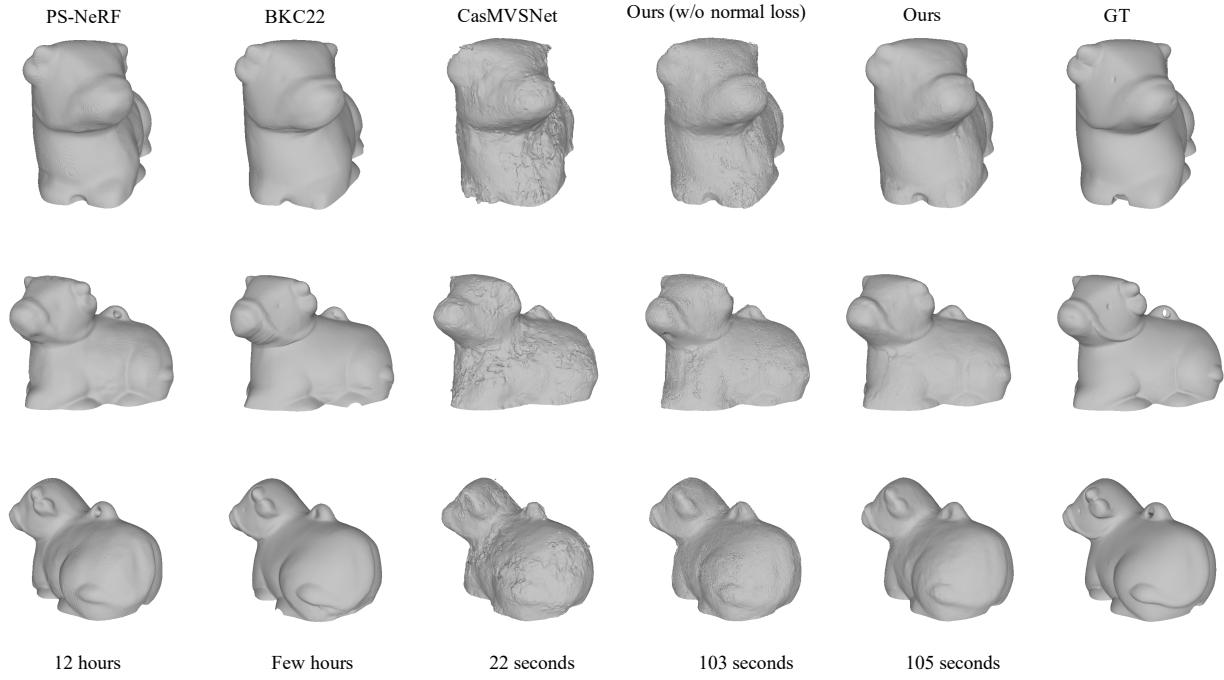


Figure 8. Reconstruction of COW under three different views (front, right-side, back) in DiLiGenT-MV [38]. Last row is reconstruction time.

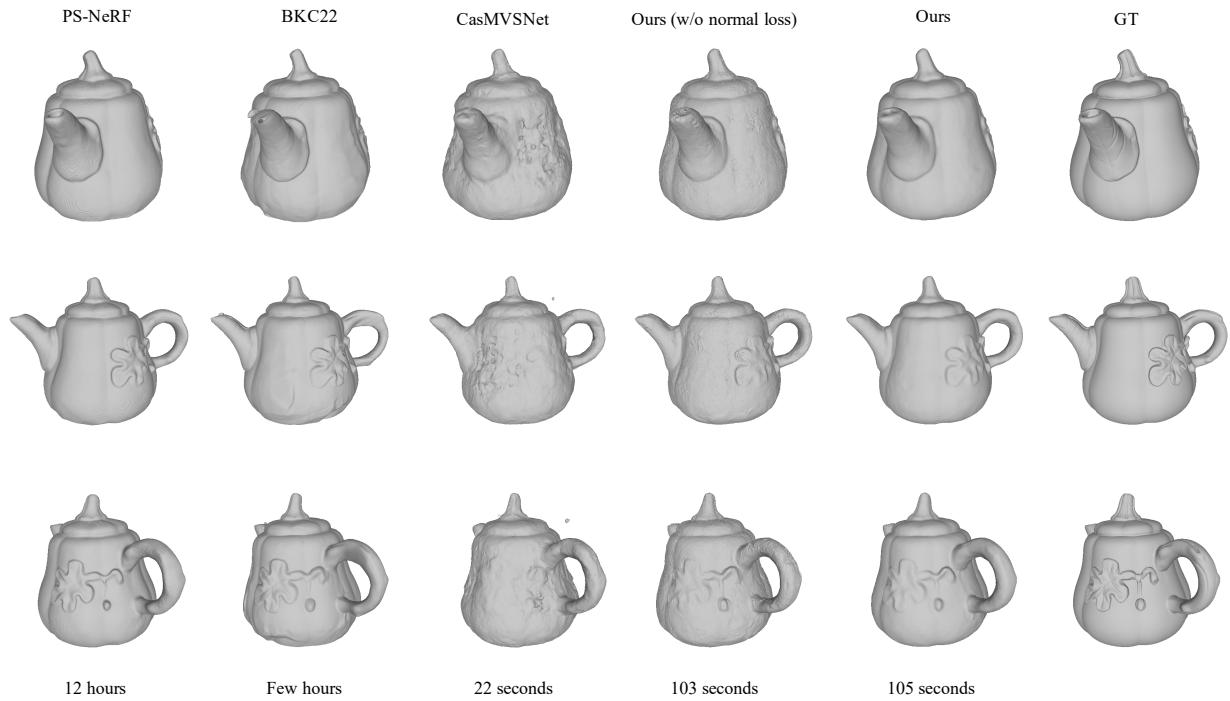


Figure 9. Reconstruction of POT2 under three different views (left-side, front, right-side) in DiLiGenT-MV [38]. Last row is reconstruction time.

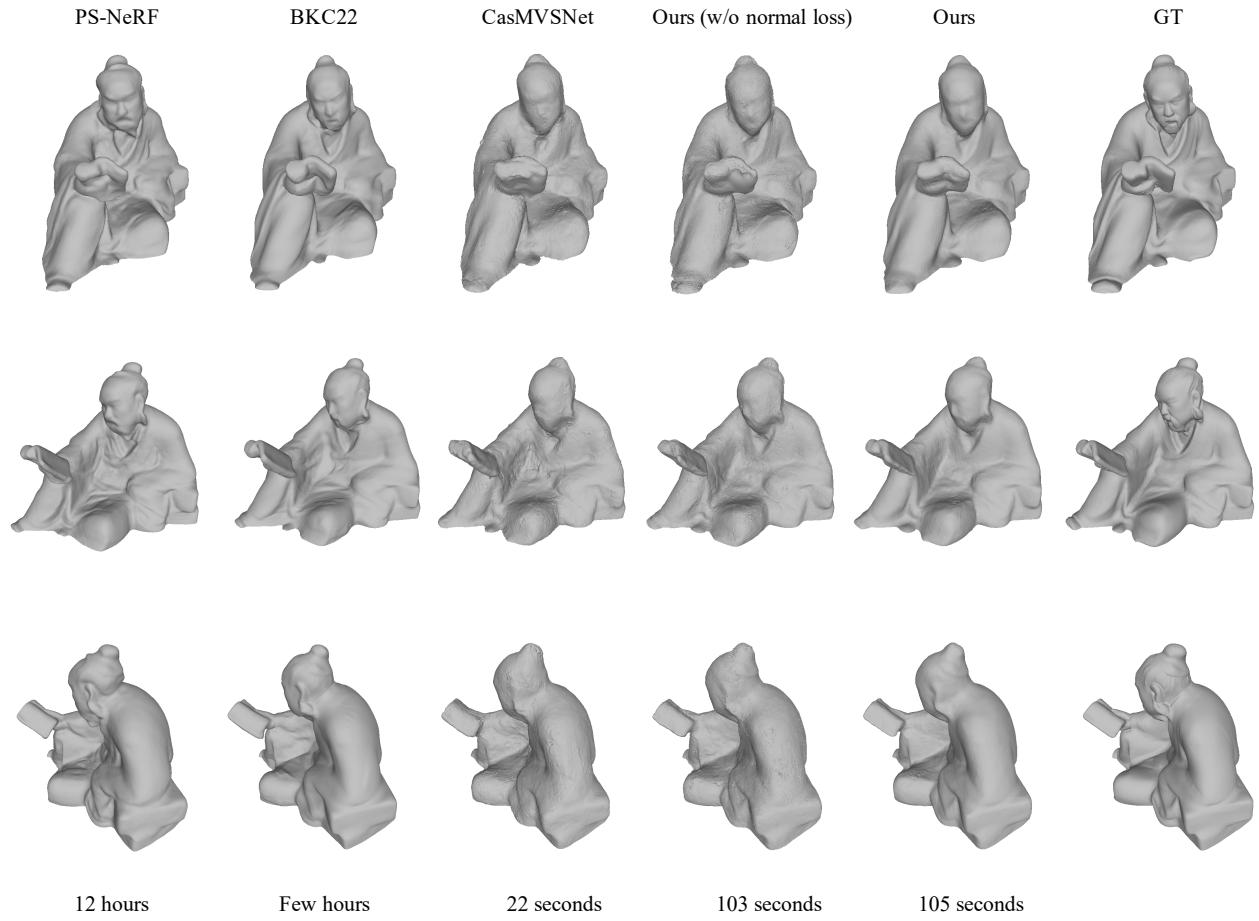


Figure 10. Reconstruction of READING under three different views (front, right-side, back) in DiLiGenT-MV [38]. Last row is reconstruction time.

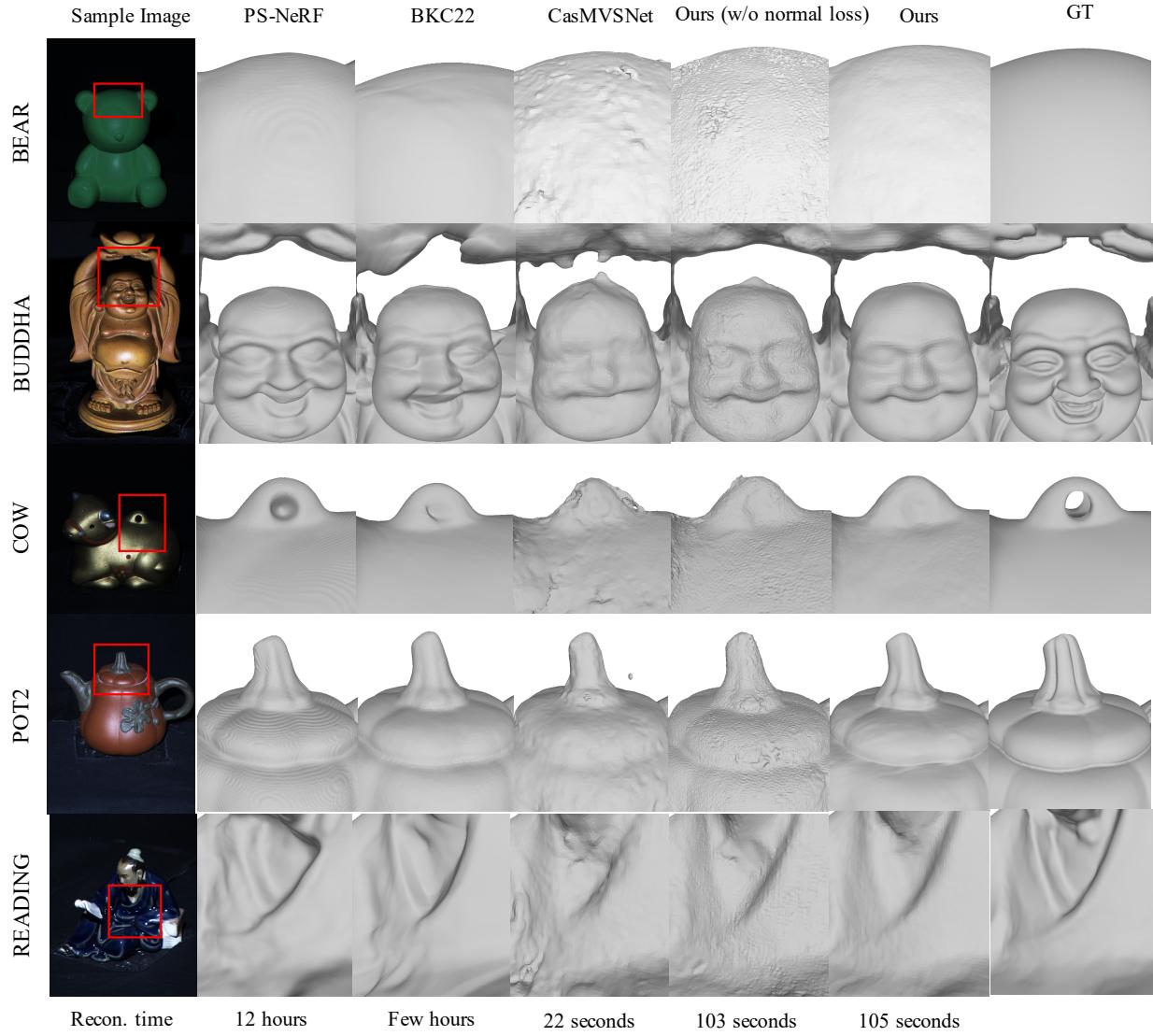


Figure 11. Zoomed-in areas on meshes from all methods. We observe that in general PS-NeRF [65] provides mesh with global fine details while it also contains iso-contour pattern artifacts. Our method can provide smooth mesh with correct global shape even though it takes very short time compared with other methods.