

# Enhanced Stable View Synthesis

Nishant Jain\*  
Indian Institute of Technology  
Roorkee, India  
njain@cs.iitr.ac.in

Suryansh Kumar\*<sup>†</sup>  
ETH Zürich  
Switzerland  
sukumar@ethz.ch

Luc Van Gool  
ETH Zürich  
Switzerland  
vangool@ethz.ch

## Abstract

We introduce an approach to enhance the novel view synthesis from images taken from a freely moving camera. The introduced approach focuses on outdoor scenes where recovering accurate geometric scaffold and camera pose is challenging, leading to inferior results using the state-of-the-art stable view synthesis (SVS) method. SVS and related methods fail for outdoor scenes primarily due to (i) over-relying on the multiview stereo (MVS) for geometric scaffold recovery and (ii) assuming COLMAP computed camera poses as the best possible estimates, despite it being well-studied that MVS 3D reconstruction accuracy is limited to scene disparity and camera-pose accuracy is sensitive to key-point correspondence selection. This work proposes a principled way to enhance novel view synthesis solutions drawing inspiration from the basics of multiple view geometry. By leveraging the complementary behavior of MVS and monocular depth, we arrive at a better scene depth per view for nearby and far points, respectively. Moreover, our approach jointly refines camera poses with image-based rendering via multiple rotation averaging graph optimization. The recovered scene depth and the camera-pose help better view-dependent on-surface feature aggregation of the entire scene. Extensive evaluation of our approach on the popular benchmark dataset, such as Tanks and Temples, shows substantial improvement in view synthesis results compared to the prior art. For instance, our method shows **1.5 dB** of PSNR improvement on the Tank and Temples. Similar statistics are observed when tested on other benchmark datasets such as FVS, Mip-NeRF 360, and DTU.

## 1. Introduction

Image-based rendering, popularly re-branded as view synthesis, is a long-standing problem in computer vision and graphics [45, 47]. This problem aims to develop a

method that allows the user to seamlessly explore the scene via rendering of the scene from a sparse set of captured images [2, 23, 45]. Furthermore, the rendered images must be as realistic as possible for a better user experience [37–39]. Currently, among the existing approaches, Riegler and Koltun stable view synthesis (SVS) approach [39] has shown excellent results and demonstrated photorealism in novel view synthesis, without using synthetic gaming engine 3D data, unlike [37]. SVS is indeed stable in rendering photorealistic images from novel viewpoints for large-scale scenes. Yet, it assumes MVS [42, 43] based dense 3D scene reconstruction and camera poses from COLMAP [42] are correct. The off-the-shelf algorithms used for 3D data acquisition and camera poses from images are, of course, popular, and to assume these algorithms could provide favorable 3D reconstruction and camera poses is not an outlandish assumption. Nonetheless, taking a step forward, in this paper, we argue that although choices made by SVS for obtaining geometric scaffold and camera poses in the pursuit of improving view synthesis is commendable, we can do better by making mindful use of fundamentals from multiple-view geometry [14, 15] and recent developments in deep-learning techniques for 3D computer vision problems.

To start with, we would like to emphasize that it is clearly unreasonable, especially in an outdoor setting, to assume that multi-view stereo (MVS) can provide accurate depth for all image pixels. It is natural that pixels with low disparity will not be reconstructed well using state-of-the-art MVS approaches [10, 42, 43, 49]. Even a precise selection of multiple view images with reasonable distance between them (assume good baseline for stereo) may not be helpful due to loss of common scene points visibility, foreshortening issue, etc. [46]. Such issues compel the practitioner to resort to post-processing steps for refining the MVS-based 3D geometry so that it can be helpful for rendering pipeline or neural-rendering network at train time.

Another critical component to view synthesis, which is often brushed aside in the literature is the accurate recovery of the camera poses. In neural view synthesis approaches such as [39], if the camera pose is wrong, the feature ag-

\*Equal Contribution

<sup>†</sup>Corresponding Author (k.sur46@gmail.com)



Figure 1. **Qualitative comparison.** Our result compared to the popular SVS method [39] on the M60 scene of the tanks and temples dataset [21]. It is easy to observe that our approach can better render fine details in the scene. For this scene, the PSNR values for SVS [39] and our method are 19.1 and 20.8, respectively, demonstrating improved PSNR result.

gregation corresponding surface points could be misleading, providing inferior results. Therefore, we should have camera poses as accurate as possible. Unfortunately, despite the camera pose importance to this problem, discussion on improving camera pose is often ignored under the assumption that COLMAP [42] provide the best possible camera pose estimates for any possible scenarios. Practically speaking, this is generally not the case for outdoor scenes [5, 14, 18]. What is more surprising is that some recent benchmark datasets put COLMAP recovered poses as the ground-truth poses [36]. Hence, we want to get this out way upfront that a robust and better camera-poses estimates are vital for better modeling view synthesis problem.

From the above predication, it is apparent that a more mindful approach is required to make view synthesis approaches practically useful, automatic, and valuable for real-world application. To this end, we propose a principled and systematic approach that provides a better geometric scaffold and camera poses for reliable feature aggregation of the scene’s surface points, leading to improved novel-view synthesis results enabling superior photorealism.

In practice, we can have suitable initial camera poses from images using COLMAP. Yet, it must be refined further for improved image-feature aggregation corresponding to 3D surface points for neural rendering. It is well-studied in multiple-view geometry literature that we can improve and refine camera poses just from image key-point correspondences [13, 14]. Accordingly, we introduce a learning-based multiple motion averaging via graph neural network for camera pose recovery, where the pose graph is initialized using COLMAP poses for refinement.

Meanwhile, it is challenging to accurately recover the 3D geometry of scene points with low or nearly-zero disparity using MVS methods [15, 46]. Another bad news from the theoretical side is that a precise estimation of scene depth from a single image is unlikely<sup>1</sup>, which is a correct statement and hard to argue. The good news is that advancements in deep-learning-based monocular depth prediction have led to some outstanding results in several practical

applications [27, 34]. Thus, at least practically, it seems possible to infer reliable monocular depth estimates up to scale. Using single image depth prediction, we can reason about the depth of scene points with low disparities. So, our proposed strategy is to use confidence based multiple-view stereo 3D that favours pixels with near-to-mid disparity and allows monocular depth estimates for the rest of the pixels. Overall depth is recovered after scaling all the scene depth appropriately using MVS reconstructed metric.

By encoding the image features via convolutional neural networks, we map the deep features to our estimated 3D geometric scaffold of the scene. Since we have better camera poses and scene reconstruction, we obtain and aggregate accurate feature vectors corresponding to each imaging view-rays—both from the camera to the surface point and from the surface point to viewing image pixels, giving us a feature tensor. We render the new image from the features tensor via a convolutional network and simultaneously refine the camera pose. In summary, our contributions are

- A systematic and principled approach for improved stable view synthesis enabling enhanced photorealism.
- The introduced approach exploits the complementary nature of MVS and monocular depth estimation to recover better 3D geometric scaffold of the scene. Meanwhile, the robust camera poses are recovered using graph neural network based multiple motion averaging.
- Our approach proposes an improved loss function to jointly optimize and refine for poses, neural image rendering, and scene representation showing superior results.

Our approach when tested on benchmark datasets such as Tank and Temples [21], FVS [38], Mip-NeRF 360 [1], and DTU [19] gives better image based rendering results with generally more than 1 dB PSNR gain (see Fig.1).

## 2. Background and Preliminaries

Our work integrates the best of novel view synthesis and multiple view geometry approaches in computer vision in a mindful way. Both novel view synthesis and MVS are classical problems in computer graphics and computer vision, with an exhaustive list of literature. Thus, we confine our re-

<sup>1</sup>As several 3D scene points can have same image projection.

lated work discussion to the methods that is directly related to the proposed approach. The interested reader may refer to [4, 10, 15, 31, 41, 45, 47] for earlier and recent progress in these areas. Here, we briefly discuss relevant methods and the current state-of-the-art in neural image-based rendering.

**(i) Uncalibrated Multi-View Stereo.** Given the intrinsic camera calibration matrix, we can recover camera poses and the 3D structure of the scene using two or more images [10, 24, 25, 42, 43]. One popular and easy-to-use MVS framework is COLMAP [42], which includes several carefully crafted modules to estimate camera poses and sparse 3D structures from images. For camera pose estimation, it uses classical image key-point-based algorithms [16, 30, 48]. As is known that such methods can provide sub-optimal solutions and may not robustly handle outliers inherent to the unstructured images. Consequently, camera poses recovered using COLMAP can be unreliable, primarily for outdoor scenes. Moreover, since the 3D reconstruction via triangulation uses sparse key points, at best accurate semi-dense 3D reconstruction of the scene could be recovered. Still, many recent state-of-the-art methods in novel view synthesis heavily rely on it [28, 38, 39, 52].

**(ii) Image-based rendering.** Earlier image-based rendering approaches enabled novel view synthesis from images without any 3D scene data under some mild assumptions about the camera, and imaging [12, 26, 44]. Later with the development of multiple-view stereo approaches and RGBD sensing modalities, several works used some form of 3D data to improve image-based rendering. Popular works along this line includes [2, 6, 17, 23, 32]. In recent years, neural network-based methods have dominated this field and enabled data-driven approach to this problem with an outstanding level of photorealism [29, 37–39, 52].

Among all the current methods we tested, SVS [39] stands out in novel view synthesis from images taken from a freely moving camera. This paper proposes an approach closely related to the SVS pipeline. SVS involves learning a scene representation from images using its dense 3D structure and camera poses predicted via uncalibrated MVS approaches detailed previously. The recovered 3D structure of the scene is post-processed to have a refined geometric scaffold in mesh representation. Using refined scene 3D, it learns a feature representation for a set of images and then projects it along the direction of the target view to be rendered. Finally, SVS involves re-projection from feature to image space resulting in the synthesized image for the target view. Given  $\mathcal{S}$ , the 3D structure for a scene, feature representation network  $f_\theta$  with  $\theta$  representing the network parameters, and the rendering network  $\mathcal{G}_\mu$  with parameters  $\mu$ , the resulting image  $\mathcal{I}_r$  along the pose  $\mathbf{p}$  is rendered as

$$\mathcal{I}_r = \mathcal{G}_\mu(\phi_a(\mathcal{S}, \mathcal{I}_s, \mathbf{p}, f_\theta)), \quad (1)$$

where,  $\mathcal{I}_s$  denotes the image set for a given seen.  $\phi_a(\cdot)$  ag-

gregates features for all the images along a given direction in  $\mathcal{I}_s$  predicted when passed through  $f_\theta$ .

Contrary to the choices made by the SVS, in this work, we put forward an approach for better estimation of the overall scene 3D geometric scaffold—for both low and high-disparity pixels—and camera poses. Our approach enables a better aggregation of surface features allowing an improved level of realism in novel view synthesis.

### 3. Method

We begin with a discussion about formulating and estimating a better scene representation by carefully exploiting the complementary behavior of monocular depth and multi-view stereo approaches in 3D data acquisition. During our implementation, we use confidence measures for precise reasoning of estimated 3D coming from both the modalities and their accuracy. Next, we detail integrating a camera pose-refining optimization to estimate better camera poses for accurately projecting our 3D scene representation to the target view. Finally, the proposed joint optimization objective to refine structure representation and camera poses simultaneously is discussed.

#### 3.1. Overview

Given a set of source images  $\mathcal{I}_s$  taken from a freely moving camera, the current state-of-the-art SVS [39] first estimates camera-poses  $\mathcal{P}$  and scene 3D structure  $\mathcal{S}$ . Next, it encodes the images via a CNN. Then, using the estimated camera, it maps the encoded features onto the scene 3D scaffold. For each point  $\mathbf{x} \in \mathbb{R}^3$  on the scaffold, SVS query set of images in which  $\mathbf{x}$  is visible to obtain its corresponding feature vectors. This feature vector is then conditioned on the output view direction via a network to produce a new feature vector. Such a new feature vector is obtained for all the points on the scaffold to form a feature tensor, which is then decoded using CNN to synthesize the output image. Denoting the complete set of parameters for learning the feature representation and rendering as  $\Theta$ , we can write SVS idea as

$$\Theta \sim \Psi(\Theta | \mathcal{S}, \mathcal{P}, \mathcal{I}_s) \circ \Psi(\mathcal{S}, \mathcal{P} | \mathcal{I}_s). \quad (2)$$

$\Psi(\cdot)$  symbolizes an abstract functional. Here,  $\mathcal{S}$  and  $\mathcal{P}$  are estimated using structure from motion and MVS [42, 43]. As mentioned before, we first aim to recover a much better 3D scene recovery by utilizing the complementary nature of the monocular depth and stereo depth. Furthermore, estimate improved camera poses using neural graph-based multiple motion averaging. We denote  $\mathcal{D}_s$  as monocular depth for each image in  $\mathcal{I}_s$  and  $\mathcal{S}_R, \mathcal{P}_R$  as our scene 3D reconstruction and camera poses, respectively.

Once we estimate  $\mathcal{S}_R, \mathcal{P}_R$ , we jointly optimize for better neural scene representation and camera poses with the neural network parameters  $\Theta$  using the proposed loss function

(Sec. 3.3). Our better scene representation and improved camera help in improved view-synthesis at test time. Our overall idea can be understood via the following equation:

$$\Theta, \mathcal{S}_R, \mathcal{P}_R \sim \Psi(\Theta, \mathcal{S}_R, \mathcal{P}_R | \mathcal{S}, \mathcal{P}, \mathcal{D}_s, \mathcal{I}_s) \quad (3)$$

### 3.1.1 3D Scene Features

Here, we describe feature aggregation for the output view obtained using the SVS pipeline [39] combined with the back-projected RGB-D features due to images and their monocular depth via a convolutional neural network (CNN). Given the output direction  $\mathbf{u}$ , for each  $\mathbf{x} \in \mathbb{R}^3$  on the geometric scaffold along  $\mathbf{u}$  there is a subset of source images, in which  $\mathbf{x}$  is visible. Assuming the total number of images in this subset to be  $K$ , we denote the set of view directions corresponding to these images as  $\{\mathbf{v}_k\}_{k=1}^K$ .

**(i) SVS Features.** The source images in  $\mathcal{I}_s$  are first passed through a feature extraction CNN to obtain a feature tensor  $\mathcal{F}_k$  corresponding to  $k^{\text{th}}$  image. Denoting  $f_k(\mathbf{x})$  as the feature corresponding to a point  $\mathbf{x} \in \mathbb{R}^3$  in  $\mathcal{F}_k$  located at the projection (or bilinear interpolation) of  $\mathbf{x}$  on  $k^{\text{th}}$  image. To this end, SVS [39] proposed the following aggregation function ( $\phi_\alpha$ ) to compute the feature for  $\mathbf{x}$  along  $\mathbf{u}$

$$\phi_\alpha(\mathbf{u}, \{\{\mathbf{v}_k, f_k(\mathbf{x})\}\}) = \frac{1}{W} \sum_{k=1}^K \max(0, \mathbf{u}^T \mathbf{v}_k) f_k(\mathbf{x}), \quad (4)$$

where,  $W = \sum_{k=1}^K \max(0, \mathbf{u}^T \mathbf{v}_k)$  is sum of all the weights.

**(ii) Monocular Features.** Given an image, we predict its depth with a per-pixel confidence score. For this, we use an existing monocular depth prediction network [35] pre-trained on Omnidata [7]. To predict a normalized depth prediction confidence score per pixel  $w_{ij}$ , we add another network on top of it (refer supplementary for details), which takes both images and depth as input.

$$\sum_i \sum_j w_{ij} \cdot \mathcal{L}(d_{ij}, \hat{d}_{ij}); \text{ where, } w_{ij} = \phi_w(\mathcal{I}_s, \mathcal{D}_s)_{ij}. \quad (5)$$

$\mathcal{L}$  symbolizes the  $l_2$  loss between the known pixel depth  $d_{ij}$  and the predicted depth  $\hat{d}_{ij}$  for  $(i, j)$  pixel in the 2D depth image. We call network  $\phi_w$  as the ‘‘confidence-prediction head’’. We take the depths predicted by this network and fuse them with the source image using the Channel Exchange Network [50]. The fused information is then projected into  $N$  dimensional feature space using network  $\mathcal{F}_\theta$ , which is a CNN-based architecture. To compute the feature corresponding to a pixel  $\mathbf{p}$  in the target view  $\mathbf{u}$ , we warp the features predicted for each source views  $\mathbf{v}_k$  to this target view but now using the monocular depth  $d \in \mathcal{D}_s$  of the source view.

$$g_k(\mathbf{u}, \mathbf{p}) = f_k^m(\mathcal{W}(\mathbf{p}, \mathbf{u}, \mathbf{v}_k, d)) \quad (6)$$

where,  $g_k$  denotes the feature warping function,  $f_k^m(p)$  corresponds to feature in the tensor  $\mathcal{F}_\theta(\mathcal{I}_s^k, d)$  (corresponding to  $k^{\text{th}}$  source image and corresponding monocular depth) at pixel  $\mathbf{p}$ , and  $\mathcal{W}$  is the warping function. We now aggregate the warped features corresponding to each source image using a weighted sum based on confidence scores:

$$\phi_\alpha(\mathbf{u}, \{\{\mathbf{v}_k, f_k^m(\mathbf{p})\}\}) = \sum_{k=1}^K c_k(\mathbf{p}) f_k^m(\mathbf{p}) \quad (7)$$

here,  $c_k$  symbolizes the predicted depth confidence map in the  $k^{\text{th}}$  view.

### (iii) Structure Feature Aggregation Per Scene Point.

Given we have features from two different modalities *i.e.*, monocular depth and MVS, we now propose an aggregated feature representation for the target view per scene point  $h_f = h_\theta(h_m, h_s)$ , where  $h_\theta$  is a CNN-based neural network with parameters  $\theta$  and  $h_m$  is the final representation for the monocular estimation and  $h_s$  for the stereo-estimation. We aim to attain maximal correlation from both of the input representations *i.e.*,

$$\lambda_1 \mathcal{C}(h_f, h_m) \cdot \mathcal{C}(h_f, h_s) + \lambda_2 \mathcal{C}(h_f, h_s) + \lambda_3 \mathcal{C}(h_f, h_m) \quad (8)$$

where,  $\sum_{i=1}^3 \lambda_i = 1$ ;  $\lambda_i > 0 \forall i \in \{1, 2, 3\}$  and  $\mathcal{C}(\cdot)$  being the standard correlation function [3]. We set  $\lambda_1 = 1$  for nearby points where confidence of monocular depth could be greater than a defined threshold ( $\tau$ ).  $\lambda_2 = 1$  for nearby points where confidence of monocular depth is lower than a defined threshold and finally  $\lambda_3 = 1$  for far points whose relative depth is greater than a pre-defined threshold ( $\sigma$ ). Such a choice is made since stereo features might be less accurate in certain depth range due to low disparity or too close to the lens. The final aggregation  $h_f$  comprises of CNN networks  $\phi_\alpha$  and  $\phi_\beta$ , parameterized by  $\alpha$  and  $\beta$ , respectively. The transformed features from each of these modalities are then fed to a CNN network  $\phi_\nu$  with parameters  $\nu$ . Overall aggregated features per scene point is represented as

$$h_f = h_\theta(h_m, h_s) = \phi_\nu(\phi_\alpha(h_m), \phi_\beta(h_s)). \quad (9)$$

### 3.1.2 Camera Pose Estimation

Given we have initially estimated pose-set  $\mathcal{P}$ , we use multiple motion averaging (MRA) to recover better camera poses. MRA is fast, efficient, and robust in camera pose estimation from noisy key-point correspondences and works well even for sparse image sets. MRA takes a global view of the camera recovery problem. A view-graph corresponding to the initial set of camera poses is constructed as a global representation, with each graph node as a pose. For pose optimization, MRA first performs rotation averaging to recover rotation and then solves for translation. Assume a

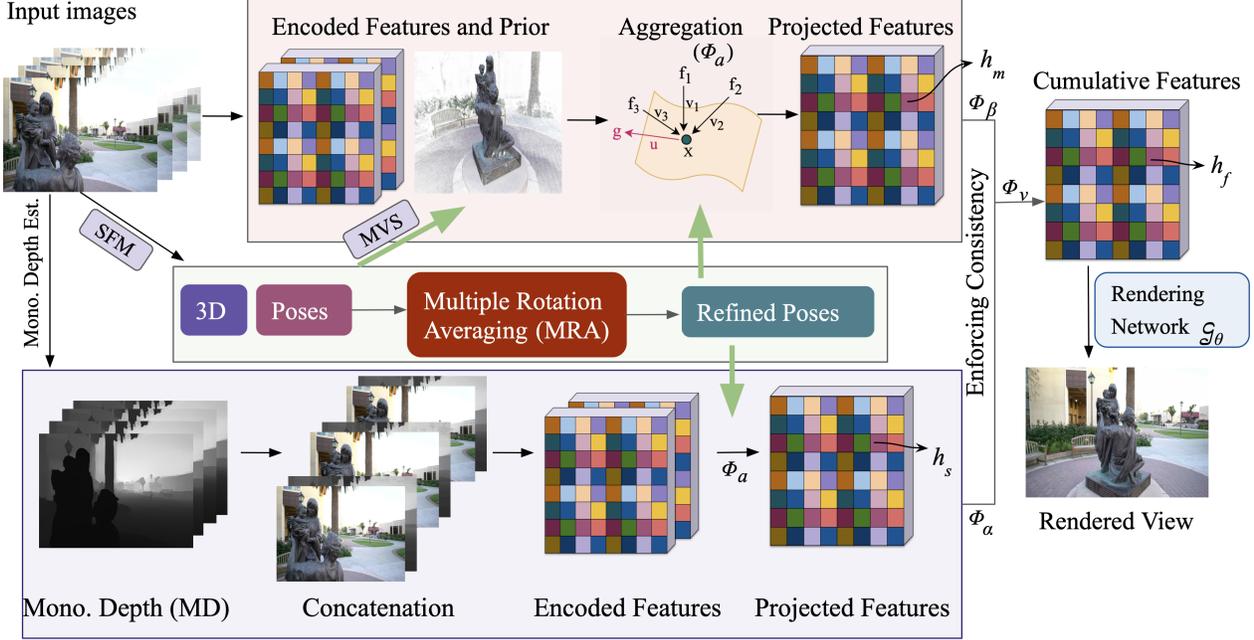


Figure 2. **Overview.** For a given set of input images, we first estimate the scene 3D structure and the initial value of camera poses using structure from motion (SFM) followed by MVS. We also estimate per-image depth using an existing Monocular Depth estimation model [7]. Next, we generate two sets of Projected Features for each target view: The first set (upper stream) by encoding the input images and then unprojecting the features onto the scaffold, followed by re-projection to the target view and aggregation. For the second set (lower stream), the images are first concatenated with their monocular depths (MD), and the encoded features are then projected using these depths, followed by aggregation. These two feature sets are then merged into Cumulative Features, and finally, the target view is rendered.

directed view-graph  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ . A vertex  $\mathcal{V}_j \in \mathcal{V}$  in this view graph corresponds to  $j^{\text{th}}$  camera absolute rotation  $R_j$  and  $\mathcal{E}_{ij} \in \mathcal{E}$  corresponds to the relative orientation  $\tilde{R}_{ij}$  between view  $i$  and  $j$ . For our problem, relative orientations are given (noisy), and we aim to solve for absolute pose  $R_j$  robustly. Conventionally, in the presence of noise, the solution is obtained by solving the following optimization problem to satisfy compatibility criteria.

$$\operatorname{argmin}_{\{R_j\}} \sum_{\mathcal{E}_{ij} \in \mathcal{E}} \rho\left(\gamma(\tilde{R}_{ij}, R_j R_i^{-1})\right) \quad (10)$$

where,  $\gamma(\cdot)$  denotes a metric on  $SO(3)$  to measure distance between two rotation and  $\rho(\cdot)$  is a robust  $l_1$  norm. We use a graph-neural network to predict the robustified pose view-graph of the scene following the camera pose refining approach of NeuRoRA [33]. First, the camera poses are checked for outliers via a cyclic consistency in the pose-graph [14]. Then, a view graph is initialized based on these filtered poses. This view graph is then optimized using the camera pose-refining network of NeuRoRA [33].

### 3.2. Rendering novel views

The aggregated feature tensor  $\mathcal{F}^a$  along a direction  $\mathbf{u}$  comprises of the features  $h_f$  (Eq. 9) for each pixel along  $\mathbf{u}$ ,

obtained using 3D scene features (discussed in Sec. 3.1.1). This tensor is now projected to the image space function using  $\mathcal{G}_\theta$ , a CNN-based network with parameters  $\theta$ . It is quite possible that the regions in the target image may not be encountered in any of the source images. For those points, the values in feature tensor  $\mathcal{F}$  are set to  $\mathbf{0}$ , and thus, they require some inpainting. Denoting the set of test images by  $\mathcal{I}_t$ , the rendered image corresponding to the  $k^{\text{th}}$  image in the test set,  $\hat{\mathcal{I}}_t^k$ , is predicted based on the following equation:

$$\hat{\mathcal{I}}_t^k = \mathcal{G}_\theta(\mathcal{F}_k^a) \quad (11)$$

Similar to stable view synthesis, we parameterize the  $\mathcal{G}_\theta$  function using a U-Net [40] style model which also deals with inpainting/hallucinating the newly discovered regions in the test image. Fig.(2) shows the complete overview of the proposed method.

### 3.3. Joint Optimization

Given  $\mathcal{I}_s$ , we train the model using our overall loss function  $\mathcal{L}$  comprising of two terms.

$$\mathcal{L} = \mathcal{L}_s + \mathcal{L}_p \quad (12)$$

(i) The **first loss term**  $\mathcal{L}_s$  encourages the network to learn better features corresponding to the scene point. The  $\mathcal{L}_s$

	Truck			M60			Playground			Train		
	PSNR↑	LPIPS↓	SSIM↑									
NeRF++ [52]	21.8	0.31	0.81	17.6	0.43	0.73	21.9	0.39	0.79	17.6	0.48	0.68
FVS [38]	21.9	0.14	0.84	15.8	0.32	0.77	21.7	0.21	0.83	17.3	0.28	0.75
SC-NeRF [20]	22.3	0.29	0.82	18.4	0.40	0.76	22.4	0.35	0.83	18.2	0.42	0.73
Point-NeRF [51]	22.7	0.14	0.87	19.6	0.21	0.85	22.2	0.19	0.83	18.6	0.16	0.83
SVS [39]	22.9	0.12	0.88	19.1	0.22	0.83	22.9	0.17	0.86	17.9	0.19	0.81
Ours	<b>24.1</b>	<b>0.12</b>	<b>0.90</b>	<b>20.8</b>	<b>0.20</b>	<b>0.89</b>	<b>23.9</b>	<b>0.14</b>	<b>0.90</b>	<b>20.1</b>	<b>0.13</b>	<b>0.88</b>

Table 1. Rendered image quality comparison with current state-of-the-art methods in novel view synthesis on the popular Tanks and Temples dataset [21]. We use the popular metrics i.e., PSNR, LPIPS and SSIM for the comparison.

	Bike		Flowers		Pirate		Digger		Sandbox		Soccertable	
	SSIM↑	LPIPS↓										
NeRF++ [52]	0.71	0.27	0.80	0.31	0.71	0.43	0.65	0.35	0.84	0.24	0.87	0.21
FVS [38]	0.61	0.28	0.79	0.27	0.69	0.37	0.68	0.24	0.78	0.32	0.82	0.21
SVS [39]	0.74	0.22	0.84	0.21	0.75	0.32	0.77	0.18	0.85	0.20	0.91	0.15
Ours	<b>0.79</b>	<b>0.19</b>	<b>0.88</b>	<b>0.18</b>	<b>0.80</b>	<b>0.29</b>	<b>0.81</b>	<b>0.16</b>	<b>0.91</b>	<b>0.17</b>	<b>0.94</b>	<b>0.13</b>

Table 2. Quantitative comparison on FVS dataset [38]. We use the popular metrics, *i.e.*, PSNR, LPIPS and SSIM for the comparison.

comprises to two objective function  $\mathcal{L}_{rgb}$  and  $\mathcal{L}_{corr}$ . The function  $\mathcal{L}_{rgb}$  measures the discrepancies between the rendering of an image in the source set and the actual image available in the set. This objective is used to train the structure estimation and rendering network parameters jointly. Whereas  $\mathcal{L}_{corr}$  is used for maximizing the correlation objective discussed in Eq.(8) to arrive at the optimal aggregated scene representation.

$$\mathcal{L}_s = \mathcal{L}_{rgb} + \mathcal{L}_{corr} \quad (13)$$

$\mathcal{L}_{corr}$  takes the negative of the objective defined in Eq.(8). The structure network parameters corresponding to functions  $f_k$  and  $f_k^m$  (cf. Sec. 3.1.1) and is updated using  $\mathcal{L}_{corr}$  and  $\mathcal{L}_{rgb}$ . The rendering network parameters corresponding to  $\mathcal{G}_\theta$  updated using only  $\mathcal{L}_{rgb}$ .

(ii) The **second loss term**  $\mathcal{L}_p$  corresponds to the camera pose refinement. We used the following loss  $\mathcal{L}_p$  to improve the camera pose estimation.

$$\mathcal{L}_p = \mathcal{L}_{mra} \quad (14)$$

As is known, the multiple motion averaging and the color rendering cost functions are directly impacted by the camera pose parameters. And therefore, the  $\mathcal{L}_{rgb}$  is inherently used as an intermediate objective in which the camera pose term is constant between the two different representations.

## 4. Experiments, Results and Ablations

Here, we discuss our experimental results and their comparison to relevant baselines. Later, we provide critical ablation analysis, a study of our camera pose estimation idea, and 3D structure recovery. Finally, the section concludes with a discussion of a few extensions of our approach<sup>2</sup>.

<sup>2</sup>Please refer supplementary for the train-test setup, implementation details regarding hyperparameters, and architectures to reproduce our results.

**Baselines.** We compare our approach with the recently proposed novel view synthesis methods that work well in practice for real-world scenes. Our baseline list includes NeRF++ [52], SC-NeRF [20], PointNeRF [51], FVS [38] and SVS [39]. All baseline methods are trained on the source image set of the testing sequence. Our method’s train and test set are the same as SVS [39].

### 4.1. Dataset and Results

For evaluation, we used popular benchmark datasets comprising real-world and synthetic datasets. Namely, we used the Tanks and Temples [21], FVS [38], Mip-NeRF 360 [1], and the DTU [19] dataset. The first two datasets contain real-world scenes, whereas the last is an object-centric synthetic benchmark dataset. Although our approach mainly targets realistic scenes, for completeness, we performed and tabulated results on a synthetic object-based DTU dataset (see supplementary). Next, we discuss datasets and results.

(i) **Tanks and Temples Dataset.** It consists of images of large-scale real-world scenes taken from a freely moving camera, consisting of indoor and outdoor scenes. Unfortunately, we have no ground-truth poses; thus, estimated poses using COLMAP [42] are treated as pseudo-ground-truth poses. The dataset consists of a total of 21 scenes. Similar to the SVS setup, we use 15 of 21 scenes for training, 2 for validation, and the rest as test scenes. We also use a disjoint set of target views from the test scenes to evaluate all the methods. The source views in test scenes are used for scene-specific fine-tuning if required. Table 1 shows the statistical comparison of our approach with competing baselines. Our approach consistently outperforms the relevant baselines on this challenging dataset. Furthermore, if we zoom into the scene’s finer details, our approach clearly preserves details better than other approaches (see Fig. 3).

(ii) **FVS Dataset.** It comprises six real-world scenes. Each

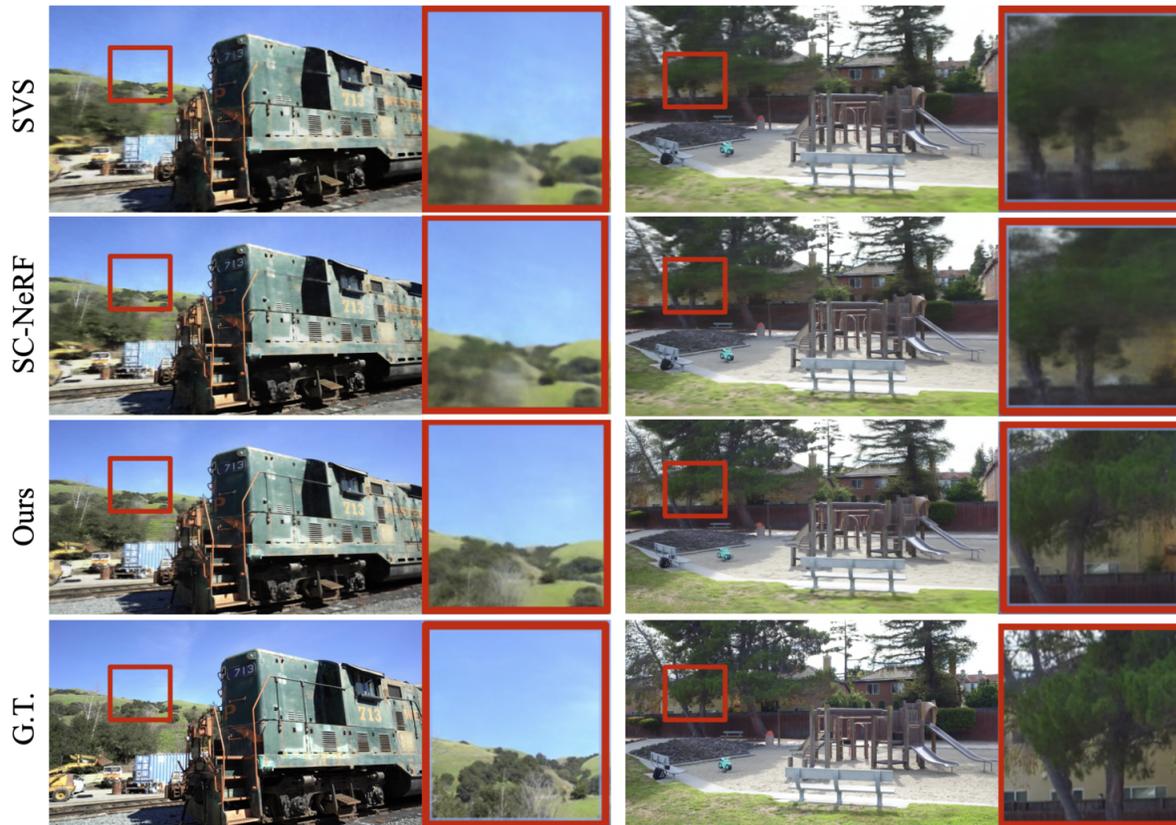


Figure 3. Qualitative comparison on tanks and temples dataset. If we zoom into the scene details, our approach results show considerably less artifacts than the state-of-the-art methods enabling unparalleled level of realism in image-based rendering. Our PSNR values for the above two scenes are (20.1, 23.9). In comparison, SC-NeRF [20] and SVS [39] provide (17.3, 22.4) and (17.9, 22.9), respectively.

scene was recorded two or more times to completely differentiate the source image set to form the target image set. Table 2 summarises this dataset’s view synthesis quantitative results. We compare our method with the state-of-the-art NeRF++ [52], Free-View synthesis [38], and SVS [39] methods. We observe that FVS improves over NeRF++, whereas SVS results on this dataset do not consistently outperform other previous methods. On the contrary, our method consistently provides better results than prior art in all categories. The dataset’s sample images and its qualitative results are provided in the supplementary.

(iii) **Mip-NeRF 360 Dataset.** We further test our approach on this recently proposed dataset comprising unbounded scenes [1]. We use the train-test setup per scene proposed by the authors to fine-tune our network and compare it with Mip-NeRF 360 [1], NeRF++ [52], and SVS [39] baseline. Table 3 shows the results for this comparison. Clearly, our methodical take on this problem generalizes well, outperforming other baseline methods for most examples.

	PSNR $\uparrow$	LPIPS $\downarrow$	SSIM $\uparrow$
NeRF++ [52]	25.03	0.355	0.682
SVS [39]	25.28	0.218	0.783
PBNR [22]	23.55	0.262	0.722
Mip-NeRF 360 [1]	27.07	0.251	0.781
Ours	<b>27.95</b>	<b>0.232</b>	<b>0.797</b>

Table 3. Quantitative comparison of our approach with existing baselines on the Mip-NeRF 360 dataset.

## 4.2. Ablation Analysis

(i) **Effect of our camera-pose estimation approach.** We performed this experiment to show the benefit of using multiple motion averaging (MRA) instead of relying only on COLMAP poses and taking them for granted. For this experiment, we used the popular Lego dataset [28]. We first compute the COLMAP camera poses shown in Fig.4 (Left). Next, we use this initial camera pose set to perform MRA utilizing the view-graph optimization. Our recovered camera poses are shown in Fig.4 (Right). For clarity, we also show the ground-truth camera pose frustum and the error between the ground-truth pose and recovered pose for the

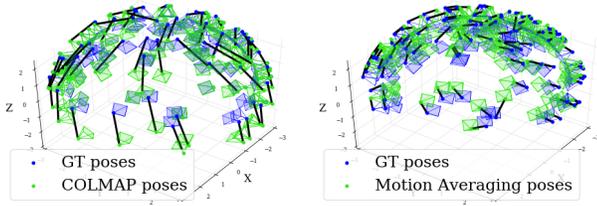


Figure 4. COLMAP [42] and our method’s camera poses on Lego dataset [28]. Here, we add some noise to the pairwise matched correspondences to simulate realistic scenarios. **Left:** Estimated camera pose from COLMAP and the ground truth pose. Here, black line shows the pose error. **Right:** Our recovered camera pose. It is easy to infer our approach robustness.

two respective cases. It is clear from Fig.(4) that a principled approach is required for better camera-pose recovery in view-synthesis problem. Furthermore, to show the benefit of our estimated camera pose in novel-view synthesis, we show our network training loss curve with and without our camera pose estimation approach in Fig.(5). For clarity, we also plot the SVS training curve. Firstly, our method shows better loss response due to the improved 3D geometric scaffold of the scene. Moreover, we improve on training loss further by introducing neural graph-based MRA, showing the benefit of our proposed loss function.

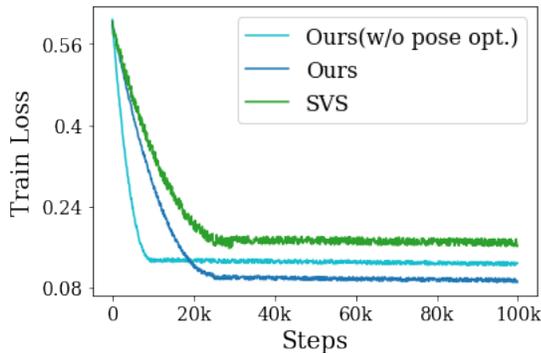


Figure 5. Training loss curve of our method with and without pose optimization compared to state-of-the-art SVS [39]

**(ii) Effect of using MVS with monocular depth.** We performed this ablation to demonstrate the effectiveness of our approach in recovering geometric reconstruction of the scene. Fig.(6) clearly shows the complementary nature of the monocular depth and stereo depth in 3D reconstruction from images. While MVS provides reliable depth results for near and mid-range points, monocular depth can deliver a good depth inference (up to scale) on the far points in the scene. By carefully integrating both modalities’ depth, we have better 3D information about the scene. Fig.(6) shows our results on a couple of scenes from Tanks and Temples dataset, demonstrating the suitability of our approach.

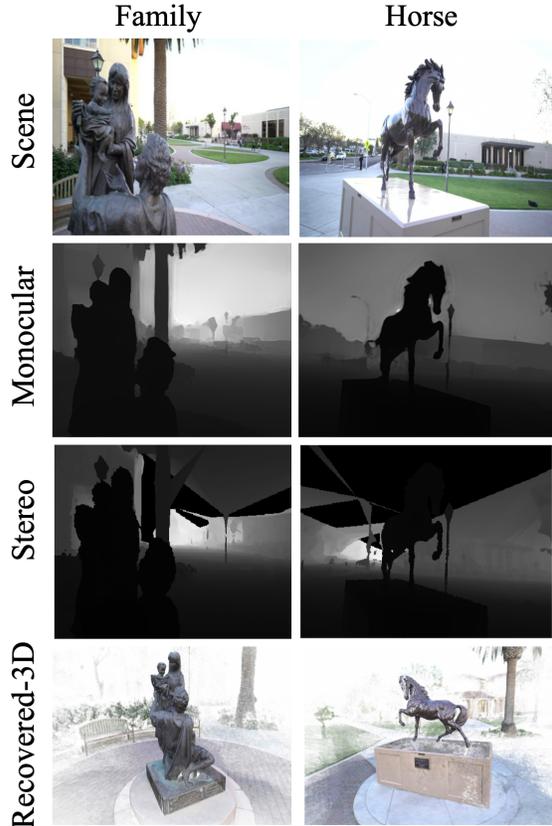


Figure 6. Comparison of depth predicted by monocular and stereo based method on the scenes from Tanks and Temples dataset.

## 5. Conclusion

The approach presented in this paper directly addresses the shortcomings of the currently popular methods in novel view synthesis. Our approach integrates concepts from the learning-based approaches and classical techniques in multiple-view geometry for solving novel-view synthesis. We demonstrate that by exploiting the (i) complementary nature of monocular depth estimation and multiple view stereo (MVS) in the 3D reconstruction of the scene, and (ii) usefulness of multiple rotation averaging in structure from motion, we can achieve better-rendering quality than other dominant approaches in novel view synthesis. We confirm the effectiveness of our approach via a comprehensive evaluation and quantitative comparison with baselines on several benchmark datasets. Although our work enabled improved photorealistic rendering, several exciting avenues exist for further improvement. One interesting future extension is fully automating our introduced view synthesis pipeline, *i.e.*, adding a learnable MVS framework to estimate intrinsic, extrinsic camera parameters with scene reconstruction for better novel view synthesis.

## References

- [1] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022.
- [2] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 425–432, 2001.
- [3] Sarath Chandar, Mitesh M Khapra, Hugo Larochelle, and Balaraman Ravindran. Correlational neural networks. *Neural computation*, 28(2):257–285, 2016.
- [4] Yuan Chang and WANG Guo-Ping. A review on image-based rendering. *Virtual Reality & Intelligent Hardware*, 1(1):39–54, 2019.
- [5] Avishek Chatterjee and Venu Madhav Govindu. Robust relative rotation averaging. *IEEE transactions on pattern analysis and machine intelligence*, 40(4):958–972, 2017.
- [6] Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM Transactions on Graphics (TOG)*, 32(3):1–12, 2013.
- [7] Ainaz Eftekhari, Alexander Sax, Jitendra Malik, and Amir Zamir. Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10786–10796, 2021.
- [8] Sun et al. Learning robust image-based rendering on sparse scene geometry via depth completion. In *CVPR*, pages 7813–7823, 2022.
- [9] Wang et al. Ibrnet: Learning multi-view image-based rendering. In *CVPR*, pages 4690–4699, 2021.
- [10] Yasutaka Furukawa, Carlos Hernández, et al. Multi-view stereo: A tutorial. *Foundations and Trends® in Computer Graphics and Vision*, 9(1-2):1–148, 2015.
- [11] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural message passing for quantum chemistry. In *International conference on machine learning*, pages 1263–1272. PMLR, 2017.
- [12] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 43–54, 1996.
- [13] Venu Madhav Govindu. Combining two-view constraints for motion estimation. In *CVPR*, volume 2. IEEE, 2001.
- [14] Richard Hartley, Jochen Trunpf, Yuchao Dai, and Hongdong Li. Rotation averaging. *International journal of computer vision*, 103(3):267–305, 2013.
- [15] Richard Hartley and Andrew Zisserman. *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [16] Richard I Hartley. In defense of the eight-point algorithm. *IEEE Transactions on pattern analysis and machine intelligence*, 19(6):580–593, 1997.
- [17] Peter Hedman, Tobias Ritschel, George Drettakis, and Gabriel Brostow. Scalable inside-out image-based rendering. *ACM Transactions on Graphics (TOG)*, 35(6):1–11, 2016.
- [18] Nishant Jain, Suryansh Kumar, and Luc Van Gool. Robustifying the multi-scale representation of neural radiance fields. In *33rd British Machine Vision Conference 2022, BMVC 2022, London, UK, November 21-24, 2022*. BMVA Press, 2022.
- [19] Rasmus Jensen, Anders Dahl, George Vogiatzis, Engin Tola, and Henrik Aanæs. Large scale multi-view stereopsis evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 406–413, 2014.
- [20] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5846–5854, 2021.
- [21] 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.
- [22] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. In *Computer Graphics Forum*, volume 40, pages 29–43. Wiley Online Library, 2021.
- [23] Johannes Kopf, Michael F Cohen, and Richard Szeliski. First-person hyper-lapse videos. *ACM Transactions on Graphics (TOG)*, 33(4):1–10, 2014.
- [24] Suryansh Kumar, Yuchao Dai, and Hongdong Li. Monocular dense 3d reconstruction of a complex dynamic scene from two perspective frames. In *Proceedings of the IEEE international conference on computer vision*, pages 4649–4657, 2017.
- [25] Suryansh Kumar, Yuchao Dai, and Hongdong Li. Superpixel soup: Monocular dense 3d reconstruction of a complex dynamic scene. *IEEE transactions on pattern analysis and machine intelligence*, 43(5):1705–1717, 2019.
- [26] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 31–42, 1996.
- [27] Ce Liu, Suryansh Kumar, Shuhang Gu, Radu Timofte, and Luc Van Gool. Va-depthnet: A variational approach to single image depth prediction. In *The Eleventh International Conference on Learning Representations (ICLR)*, 2023.
- [28] 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.
- [29] 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.
- [30] David Nistér. An efficient solution to the five-point relative pose problem. *IEEE transactions on pattern analysis and machine intelligence*, 26(6):756–770, 2004.
- [31] Onur Özyeşil, Vladislav Voroninski, Ronen Basri, and Amit Singer. A survey of structure from motion\*. *Acta Numerica*, 26:305–364, 2017.

- [32] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)*, 36(6):1–11, 2017.
- [33] Pulak Purkait, Tat-Jun Chin, and Ian Reid. Neurora: Neural robust rotation averaging. In *European Conference on Computer Vision*, pages 137–154. Springer, 2020.
- [34] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. *ICCV*, 2021.
- [35] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020.
- [36] Jeremy Reizenstein, Roman Shapovalov, Philipp Henzler, Luca Sbordone, Patrick Labatut, and David Novotny. Common objects in 3d: Large-scale learning and evaluation of real-life 3d category reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10901–10911, 2021.
- [37] Stephan R Richter, Hassan Abu Al Haija, and Vladlen Koltun. Enhancing photorealism enhancement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [38] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *European Conference on Computer Vision*, pages 623–640. Springer, 2020.
- [39] Gernot Riegler and Vladlen Koltun. Stable view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12216–12225, 2021.
- [40] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [41] Muhamad Risqi U Saputra, Andrew Markham, and Niki Trigoni. Visual slam and structure from motion in dynamic environments: A survey. *ACM Computing Surveys (CSUR)*, 51(2):1–36, 2018.
- [42] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [43] Johannes L Schönberger, Enliang Zheng, Jan-Michael Frahm, and Marc Pollefeys. Pixelwise view selection for unstructured multi-view stereo. In *European conference on computer vision*, pages 501–518. Springer, 2016.
- [44] Steven M Seitz and Charles R Dyer. View morphing. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 21–30, 1996.
- [45] Harry Shum and Sing Bing Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing 2000*, volume 4067, pages 2–13. SPIE, 2000.
- [46] Richard Szeliski. *Computer vision: algorithms and applications*. Springer Nature, 2022.
- [47] Ayush Tewari et al. State of the art on neural rendering. In *Computer Graphics Forum*, volume 39, pages 701–727. Wiley Online Library, 2020.
- [48] Bill Triggs, Philip F. McLauchlan, Richard I. Hartley, and Andrew W. Fitzgibbon. Bundle adjustment - a modern synthesis. In *Proceedings of the International Workshop on Vision Algorithms: Theory and Practice, ICCV '99*, pages 298–372, London, UK, UK, 2000. Springer-Verlag.
- [49] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. Patchmatchnet: Learned multi-view patchmatch stereo. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14194–14203, 2021.
- [50] Yikai Wang, Wenbing Huang, Fuchun Sun, Tingyang Xu, Yu Rong, and Junzhou Huang. Deep multimodal fusion by channel exchanging. *Advances in Neural Information Processing Systems*, 33:4835–4845, 2020.
- [51] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixian Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-nerf: Point-based neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5438–5448, 2022.
- [52] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020.

# Enhanced Stable View Synthesis

## —Supplementary Material—

### Abstract

*This draft accompanies the main paper. It provides more experimental results showing the suitability of our proposed approach. Furthermore, it discusses the graph neural network-based multiple rotation averaging and our software implementation details.*

## 6. Synthetic Objects

We further evaluate our proposed method on a synthetic object-centric dataset to investigate its advantages in cases where all points lie within a threshold distance from the camera. The aim is to examine whether integrating a monocular depth with the estimated stereo structure is helpful for these cases. We use the same setup as SVS [39] for evaluation, separating ten images as target novel images and use the remaining 39 as the source images for evaluating interpolation and extrapolation. Table 4 compares all the baselines and our method on this dataset. It consists of results for both view interpolation (left value) and extrapolation setups (right value) for this dataset as done in the SVS paper [39]. It can be observed that even in this case, which doesn't involve *far-away* points, our method is either similar in performance or performs better, especially in terms of PSNR values. Also, performance gap on the extrapolation task is marginally higher than the interpolation counterpart, for the PSNR values. This further highlights the importance of our method for the nearby region where we try to maximize the consistency between RGB-D features and stereo-estimated projection of image features, for places where the monocular network is highly confident.

## 7. Scene-Agnostic Model

We analyze the scene-agnostic version of our approach and SVS, also comparing with the FVS [38] method. The model is trained on scenes corresponding to training data and then is directly evaluated on a disjoint set of test scenes without any tuning. Table 5 shows the results for this version on the four scenes of the Tanks and Temples data set used in the paper, namely *truck*, *M60*, *playground* and *train*, where the model is trained using the 15 other scenes from this dataset. It can be observed that our approach can offer significantly better results, even in the scene-agnostic setup, when compared with SVS and FVS. Also, for both our method and SVS, the results are improved from scene-specific finetuning compared to the scene-agnostic setup, which can be observed by comparing the statistics presented

in Table 5 here and Table 1 in the main paper.

## 8. Graph Neural Networks for MRA

We now discuss our pose refining scheme inspired from NeuRoRA [33]. It uses a Message Passing Neural Network (MPNN) to predict robust poses given a completely initialized view graph. Given the estimated relative rotations using an SFM algorithm, they are used to initialize absolute rotations by fixing a source vertex as the frame of reference and then calculating absolute rotation of each vertex w.r.t. this frame by traversing along the minimum spanning tree. This is followed by a cyclic consistency check to remove outliers. Finally, we have the initialized observed relative rotations and initialized absolute rotations. These comprise a completely initialized view-graph.

Now, for each node  $k$  in this graph, with neighbouring set denoted by  $\mathcal{Q}_k$ , the state of this node at step  $t$ , denoted by  $h_k^t$ , is generated by processing the aggregated signal feature  $s_k^t$  it receives from all the nodes  $v \in \mathcal{Q}_k$  and its state at step  $t - 1$ :

$$h_k^t = \rho(h_k^{t-1}, s_k^t) \quad (15)$$

where  $\rho$  is some function to process these features jointly. The aggregated signal  $s_k^t$  is just a processed combination of updated states of edges corresponding to this node  $k$ :

$$s_k^t = \psi_{i \in \mathcal{Q}_k}^a \psi^b(h_i^{t-1}, h_k^{t-1}, r_{ij}) \quad (16)$$

where,  $\psi^b$  is a processing function responsible for updating the signal accumulated from the edge between nodes  $i$  and  $k$ ,  $r_{ij}$  is a feature representation for this edge. This is followed by a differentiable operation  $\psi^a$ , which can be interpreted as some activation function. For our case, both  $\rho$  and  $\psi^b$  are concatenation with 1D convolutions and a ReLU activation. Please refer [11, 33] for further details.

**Pose-refining GNN.** Given the completely initialized view graph, we denote the rotations corresponding to its vertices as  $\tilde{R}_i$ . Also, the edges of these view graph comprise the relative rotations between the 2 nodes. The edge feature  $r_{ij}$  described above is the calculated using these observed relative rotation between the two nodes denoted as  $\tilde{R}_{ij}$ . The input to the GNN is the set of rotations  $\tilde{R}_i$  and the edge feature  $r_{ij}$ . This edge feature is calculated as the discrepancy between the initialized absolute rotations  $\tilde{R}_i$  and observed relative rotations  $\tilde{R}_{ij}$  as follows:

$$r_{ij} = \tilde{R}_j^{-1} \tilde{R}_{ij} \tilde{R}_i \quad (17)$$

This leads to a supervised learning problem for the GNN, where, using the input graph denoted as  $\{\tilde{R}_i, r_{ij}\}$ , the aim

	65			106			118		
	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑	PSNR↑	LPIPS↓	SSIM↑
FVS [38]	30.1/25.2	0.03/0.07	0.97/0.95	32.3/27.1	0.03/0.08	0.95/0.93	34.9/28.7	0.02/0.07	0.97/0.94
SVS [39]	31.9/26.1	0.02/0.06	0.97/0.95	33.8/29.7	0.02/0.04	0.98/0.95	36.7/30.8	0.02/0.05	0.97/0.96
Ours	32.4/26.9	0.03/0.07	0.98/0.96	34.3/30.6	0.02/0.03	0.98/0.96	37.1/31.3	0.02/0.05	0.97/0.96

Table 4. Performance comparison with on DTU dataset [17]. We use the popular metrics i.e., PSNR, LPIPS and SSIM for the comparison.

	Truck			M60			Playground			Train		
	PSNR↑	LPIPS↓	SSIM↑									
FVS [38]	21.9	0.14	0.84	15.8	0.32	0.77	21.7	0.21	0.83	17.3	0.28	0.75
SVS [39]	22.2	0.16	0.85	18.4	0.24	0.79	22.4	0.20	0.83	17.3	0.21	0.79
Ours	<b>23.4</b>	<b>0.14</b>	<b>0.88</b>	<b>19.6</b>	<b>0.23</b>	<b>0.85</b>	<b>22.6</b>	<b>0.18</b>	<b>0.89</b>	<b>19.2</b>	<b>0.16</b>	<b>0.84</b>

Table 5. Performance comparison with state-of-the-art methods on Tanks and Temples dataset [21] in a scene-agnostic setup. We use the popular metrics i.e., PSNR, LPIPS and SSIM for the comparison.

is to estimate the absolute rotations  $\hat{R}_i$  as close as possible to the correct rotations  $\{R_i\}$  in the source node frame:

$$\{\hat{R}_i\} = \mathcal{G}(\{\tilde{R}_i\}; \Theta) \quad (18)$$

where  $\Theta$  denote the pose-GNN parameters. The network is trained for this setup using the rotation averaging loss described in the paper. Specifically, the goal is to minimize the discrepancy between observed relative rotations  $\tilde{R}_{ij}$  and estimated relative rotations  $\hat{R}_j\hat{R}_i^{-1}$ . Given only relative rotations are used in this loss function, this it might be same even if any constant angular deviation to the predicted rotations. Thus following NeuRoRA [33], we also add a weighted regularizer term to learn a one-to-one mapping between inputs and outputs which minimizes the discrepancy between initialized absolute rotations and predicted absolute rotations. This leads to the following aggregated cost function  $\mathcal{L}_{mra}$  for a given graph with  $\mathcal{E}$  denoting its edge set and  $\mathcal{V}$  denoting the set of nodes:

$$\mathcal{L}_{mra} = \sum_{\mathcal{E}_{ij} \in \mathcal{E}} d_Q(\hat{R}_j\hat{R}_i^{-1}, \tilde{R}_{ij}) + \beta \sum_{\mathcal{V}_i \in \mathcal{V}} d_Q(\hat{R}_i, \tilde{R}_i) \quad (19)$$

where  $d_Q$  is some distance metric between two rotations. We also follow quaternion representation for these rotations similar to NeuRoRA [33].

## 9. Adapting to other methods

To further show the effectiveness of our joint feature estimation and pose updating proposition, we experimented with other popular frameworks, namely IbrNet [9] and S-IbrNet [8], on two scenes of the Tanks and Temples dataset namely Truck and Playground (P.G.), in our setup. This is done by updating their feature generation modules and integrating our pose-refining module for joint optimization. Table 6 shows the PSNR values for this experiment showing numbers for these methods before and after (E-IbrNet, E-SIbrNet) the integration and also compares these results

with the method proposed in the paper. It can be observed that both the methods (IbrNet, S-IbrNet) have their PSNR values improved significantly (around 1.7, 0.9 on average across the two scenes) and finally, our proposed method in the paper performs the best on both the scenes.

	IbrNet [9]	S-IbrNet [8]	E-IbrNet	E-SIbrNet	Ours
Truck	19.7	22.5	21.9	23.6	<b>24.1</b>
P.G.	22.2	23.1	23.3	23.7	<b>23.9</b>

Table 6. PSNR comparison. E-IbrNet and E-SIbrNet show the results when our approach is put to [9] and [8] network design.

## 10. Ablation on depth threshold ( $\sigma$ )

For all the experiments, we have set the depth threshold value based on our empirical observations. We extensively analyzed this quantity on the Tanks and Temples dataset and have set it as 0.66 (relative to median depth) for all the datasets used in the paper. Table 7 shows the analysis of PSNR values on two scenes of Tanks and Temples dataset (Truck, M60) for various values of this quantity (relative to the median depth). It can be observed that the best results are at the value of 0.66, thereby providing justification for the selected value.

$\sigma$	0.25	0.50	0.66	1.0	1.25	1.5	2.0
Truck	23.4	23.8	<b>24.1</b>	23.6	23.6	23.5	23.5
M60	20.4	20.3	<b>20.8</b>	20.4	20.5	20.2	20.3

Table 7. PSNR values for ablation on the depth threshold value.

## 11. Training and Evaluation details

We now discuss the implementation details involving the training and evaluation setup along with values for parameters involved in our approach. The training is performed on the tanks and temples dataset for our method, SVS [39] and FVS [38]. Then, for the results on all the

scenes corresponding to various datasets discussed in the paper, we tune our method and SVS for scene-specific Network fine-tuning as discussed in SVS [39]. For the baselines including NeRF++ [52], SC-NeRF [20] and Point-NeRF [51] also require per-scene fitting. This scene-specific training/tuning involves using a source image set of that scene for learning and a disjoint test set of images corresponding to the same scene for evaluation. For the scene agnostic scenario, trained models on tanks and temples are directly evaluated on the test set of the scene. Also, the  $\mathcal{L}_{rgb}$  loss term, used in the Eq.(13) in the main paper, corresponds to the perceptual loss described in Eq.(6) in the SVS paper [39].

**Architecture details.** The monocular depth prediction is a DPT [35] architecture trained on Omnidata [7]. The network for predicting confidence score for each pixel is just-another head starting from fifth last layer of the depth prediction network with the same architecture as that of the depth prediction head. This is trained while keeping the depth prediction network as frozen. Note, the complete architecture is just used for obtaining depth and confidence per pixel and then is linked with rest of the pipeline. Also, the confidence weighted loss involves  $l_2$  regularization of the predicted weights to avoid the solutions, where each/most of the weights are assigned to a single/some pixels. The network  $\mathcal{F}_\theta$  for projecting monocular features consists of a ResNet-50 architecture equipped with the Channel Exchanging Layers [50]. The functions  $\phi_\alpha$ ,  $\phi_\beta$  and  $\phi_\mu$  are each 3 layered CNN networks followed by a BatchNorm layer and a ReLU activation. For the SVS features, the architecture used in the SVS paper is followed involving a U-Net for encoding images. As discussed in the paper, the rendering network is also a U-Net architecture same as in SVS paper [39].

**Hyperparameter details.** The training is performed using an Adam optimizer setting learning rate to  $10^{-3}$ ,  $\beta_1$  to 0.9,  $\beta_2$  to 0.999 and  $\epsilon$  to  $10^{-8}$ . The model is trained for 600,000 iterations on the training set with batch size of 1 and 3 source images sampled per iteration, following the SVS setup [39]. The tuning on the testing scene is carried out for 100,00 iterations. The confidence threshold parameter  $\tau$  is set to 0.05. The predefined depth threshold ( $\sigma$ ) for applying Eq.(8) in the main paper is two-thirds the median depth of the scene.