# RGBD-Net: Predicting color and depth images for novel views synthesis

Phong Nguyen
University of Oulu

Animesh Karnewar
TomTom

Lam Huynh
University of Oulu

Esa Rahtu
Tampere University

Jiri Matas
Czech Technical University in Prague

Janne Heikkilä
University of Oulu

## Abstract

*We address the problem of novel view synthesis from an unstructured set of reference images. A new method called RGBD-Net is proposed to predict the depth map and the color images at the target pose in a multi-scale manner. The reference views are warped to the target pose to obtain multi-scale plane sweep volumes, which are then passed to our first module, a hierarchical depth regression network which predicts the depth map of the novel view. Second, a depth-aware generator network refines the warped novel views and renders the final target image. These two networks can be trained with or without depth supervision. In experimental evaluation, RGBD-Net not only produces novel views with higher quality than the previous state-of-the-art methods, but also the obtained depth maps enable reconstruction of more accurate 3D point clouds than the existing multi-view stereo methods. The results indicate that RGBD-Net generalizes well to previously unseen data.*

## 1. Introduction

Novel View Synthesis (NVS), also called Image-Based Rendering (IBR), is a long-standing problem that has applications, e.g., in free-viewpoint video, telepresence, and mixed reality [47]. NVS is a problem where visual content is captured from one or several reference views and rendered from an unseen target view. The problem is challenging since mapping between views depends on the 3D geometry of the scene and the camera poses between the views. Moreover, NVS requires not only propagation of information between views but also hallucination of details in the target view that are not visible in reference image due to occlusions or limited field of view.

Early NVS methods produced target views by interpolating in ray [28] or pixel space [8]. They were followed by works that leveraged certain geometric constraints such as epipolar consistency [4] for depth-aware warping of the input views. However, these interpolation based methods suffered from artifacts arising from occlusions and inaccurate geometry. Later works tried to patch the artifacts by propagating depth values to similar pixels [5] or by soft 3D reconstruction [42].

More recently, neural rendering methods [50, 32, 34, 31] have employed deep networks to learn implicit scene representations from a large set of observations of a specific scene. Although these methods produce impressive novel views, dedicated per-scene training is required to apply the representation to a new scene. Another research direction [60, 33, 10] uses a small number of observations at each training step. While the quality of the generated novel views of these methods is worse than those produced by the neural rendering methods, they generalize to unseen data without fine-tuning or retraining.

In this paper, we bridge the gap between these two approaches and develop a method that renders high-quality novel views from an unstructured set of reference images, without needing per-scene training. Experiments show that it generalizes well to arbitrary scenes. We propose a new method called RGBD-Net that produces both color (RGB) and depth (D) images of the unseen target view. As illustrated in Fig. 1, RGBD-Net includes two main modules: a depth regression network $R$ and a depth-aware generator network $G$. The first module estimates the target view depth map and the second module refines the warped image to produce photorealistic novel views. Because of its exceptional ability of generating RGB-D images from arbitrary viewpoints, RGBD-Net also provides a new approach for 3D reconstruction. To summarize, the main contributions of our work are:

- A coarse-to-fine target-view depth regression network $R$ based on the plane sweep volume representation with a novel adaptive depth plane resampling.

- A novel depth-aware generator network $G$ utilizing the estimated depths and image priors to refine the rendered output.
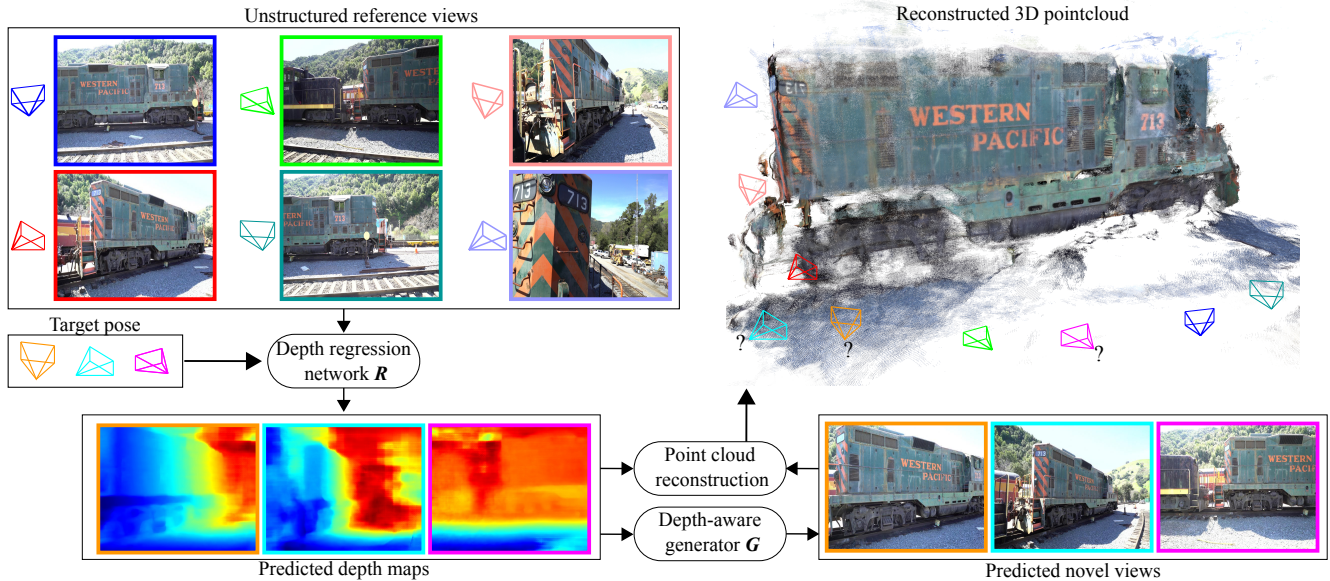
- State-of-the-art results both in novel view synthesis

Figure 1. **RGBD-Net** is a novel view synthesis method that estimates both the depth maps and color images of novel views from significantly different viewpoints. From the estimates, we reconstructed a point-cloud of the Train scene of the Tanks and Temples dataset [27]. Note: We train RGBD-Net on multi-view stereo datasets [1, 57] and the Tanks and Temples dataset is used only in evaluation.

and in multi-view 3D point cloud reconstruction.

- RGBD-Net's ability to be trained without depth supervision if ground-truth depth is not available.

Source code and neural network models will be made publicly available upon publication of the paper

## 2. Related Work

The literature related to novel view synthesis is vast and cannot be fully covered here. Interested readers are referred, e.g. to [52] for further information. We focus on reviewing method, both with and without deep learning, most relevant to our work.

***View synthesis without deep learning.*** The problem of novel view synthesis has been studied for several decades [16, 8, 28]. Many different approaches that map reference views to the novel views have been proposed. Early works on light-field based synthesis [15, 28, 46] do not require information about the scene geometry but rely on a dense and regularly-spaced camera grid. Later work by Heigl et al. [20] estimates depth maps for each input view via stereo matching and uses them for view synthesis. The common property between these methods is the restriction on the small baseline between the input views. In this paper, we focus on a less constrained set up of input views distributed around the scene.

Unstructured view synthesis often relies on 3D proxy geometry to render the novel view. Buehler et al. [4] utilizes

dense and accurate 3D geometry to map and blend input images to the target viewpoint. Later work by Chaurasia et al. [5] estimates per-view depth maps and uses these to map color values into the target view. This method also leverages superpixels to predict the missing depth values. Rather than estimating the depth map from input views, techniques such as [42, 19] use an RGB-D sensor to capture both color and depth images from different viewpoints. This paper presents a learning-based method to predict both the color and depth images of the target view using only RGB input observations.

***View synthesis with deep learning.*** Deep learning has become an essential part of image-based rendering methods to blend input views to the novel views [18, 54, 12, 3, 24] and to construct neural scene representation [2, 34, 38, 48, 50, 53, 35].

Early works on view synthesis with deep learning often use a Plane Sweep Volume (PSV) [11] for image-based rendering. Kalantari et al. [24] construct it from input views with the same viewing direction. They use an expensive light-field camera to capture these images. Each input image is projected onto successive virtual planes of the target camera to form a PSV. The mean and standard deviation per plane are utilized to estimate the disparity map and render the target view. RGBD-Net is related to this method as we also use PSVs, but we propose a hierarchical depth regression network that predicts the depth map of the novel view in a coarse-to-fine manner. Moreover, our work fo-

cuses on solving the problem of view synthesis using unstructured inputs which poses a more demanding challenge than using grid-sampled views captured by the light-field camera. Extreme View Synthesis (EVS) [10] builds upon DeepMVS [22] to estimate a depth probability volume for each input view that is then warped and fused into the target view. The initially estimated novel view is refined by comparing it with candidate patches from source images. Rather than estimating the depth maps of the source images, we train RGBD-Net to predict the depth map at the target view directly and then refine the warped novel images using a depth-aware generator network.

A significant number of works [60, 51, 55, 12] on view synthesis often represent the 3D scene by Multiple Plane Images (MPIs). Each MPI includes multiple RGB-$\alpha$ planes, where each plane is related to a certain depth. The target view is generated by using alpha composition [43] in the back-to-front order. Zhou et al. [60] introduce a deep convolutional neural network to predict MPIs that reconstructs the target views for the stereo magnification task. Srinivasan et al. [51] extend the above method and shows the relationship between the range of views that can be rendered from a multi-plane image and the depth plane sampling frequency. Later work by Flynn et al. [12] considerably improves the quality of synthesized images in the light-field setups. They present a novel network with a regularized gradient descent method to refine the generated images gradually. Local Light Field Fusion (LLFF) [33] introduces a practical high-fidelity view synthesis model that blends neighboring MPIs to the target view. The input to the MPI-based methods is also PSVs. However, those PSVs are constructed on a fixed range of depth values. The proposed RGBD-Net builds multi-scale PSVs which use adaptive sampled depth planes.

Recent geometric deep learning methods learn to deal with 3D scenes using various types of 3D representations such as voxel-grids, meshes, point-clouds, and implicit functions. These methods [32, 49, 36, 37] often use 3D convolution layers to learn 3D spatial transformations from the input views to the novel view and then perform adversarial training to enhance the quality of the output image. Another line of work [38, 31, 29, 2, 34, 35] train neural scene representations. Aliev et al. [2] recently presented Neural Point-Based Graphics (NPBG) that trains a neural network to learn feature vectors that describe 3D points in a scene. These learned features are then projected onto the target view and fed to a rendering network to produce the final novel image. The current state-of-the-art method Neural Radiance Fields (NeRF) by Mildenhall et al. [34] represents the plenoptic function by a multi-layer perceptron that can be queried using classical volume rendering to produce novel images. Despite the high quality of the synthesized novel images, these methods do not generalize well on un-

seen testing data. In contrast, RGBD-Net achieves good performance not only on a testing set that is separate from the training set but also on completely new scenes outside those datasets.

Perhaps, the closest work to RGBD-Net is the recently published Free View Synthesis (FVS) by Riegler et al. [44]. In this work, they use a structure-from-motion method [45] to reconstruct a 3D mesh of the scene for creating an incomplete depth map for the target view. They also propose a recurrent blending network to refine the warped novel views. RGBD-Net estimates a complete depth map to refine the warped novel image.

## 3. Proposed method

This section describes in detail the architecture of RGBD-Net, which comprises of two modules: a hierarchical depth regression network $R$ (Section 3.1) that estimates the depth map of the novel view and a depth-aware refinement network $G$ (Section 3.2) that enhances the warped images to produce the final target image. Last, we discuss the loss functions used to train the model in Section 3.3.

### 3.1. Depth regression network $R$

We first describe our pipeline (see Fig. 2) for estimating the depth map $\hat{D}_q$ of the target view $s_q$ from a set of unstructured input images and their poses $\{I_n, s_n\}_{n=1,...,N}$. Each reference view $I_n$ is first fed to the Feature Pyramid Network [30] to extract $K$ multi-scale features $F_n^k$ [25]. We then apply homography warping to each feature map of $F_n^k$ to construct a PSV $P_n^k$ of the target view $s_q$ with a set of $M_k$ hypothesis depth planes. Inspired by the recent work on multi-view stereo [56, 17], we estimate the depth map of the novel views in a coarse-to-fine manner. A mean PSV $\bar{P}^k = \sum_{n=1}^{N} P_n^k / N$ is fed to a 3D U-Net to estimate the novel depth map $\hat{D}_q^k$. To obtain high quality depth map, we propose a novel adaptive depth plane resampling.

**Depth plane resampling.** Previous works on view synthesis [33, 51] often construct their PSVs using a fixed number of depth planes, and they are sampled with a fixed range of depth values. In principle, we would like to sample as many hypothesis depth planes as possible to cover the whole 3D environment densely. However, the GPU memory and runtime would grow cubically as the resolution of the feature maps and the number of depth planes increases, limiting our method from generating high-resolution images. On the other hand, reducing the number of sampled planes would solve the memory issue, but it would limit the representative power of the PSV. Therefore, we propose a method to construct multi-scale PSVs using adaptive depth ranges.

At all stages, depth planes are regularly sampled with a fixed depth interval. The first stage takes the low-resolution image features to construct a PSV with a predetermined
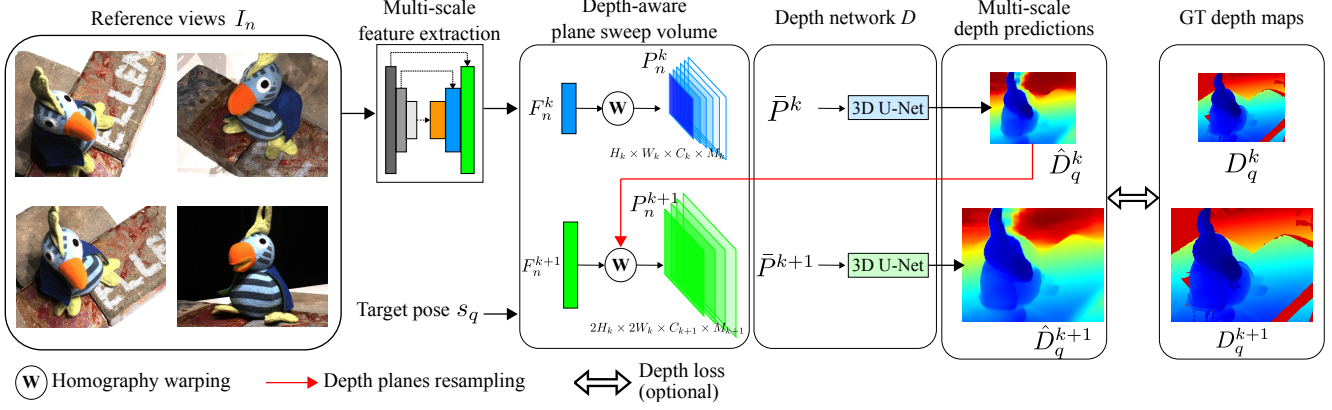
Figure 2. The architecture for estimating the depth map of the novel view. The method does not require the ground-truth depth loss for training. With only the RGB image loss, the method predicts plausible depth maps at target views.

depth range and depth interval. We then build a higher-resolution PSV using a smaller depth range with denser sampling at the following stage. More specifically, the depth planes $d_i^1$ at the scale $k = 1$ are sampled from the initial depth range as follows:

$$d_i^1 = d_{min}^1 + i\Delta_1, \quad i = 1, .., M_1 \tag{1}$$

where $d_{min}^1$ and $\Delta_1$ are the minimum depth value and depth interval, respectively. We define a sufficiently large $M_1$ and $\Delta_1$ to cover a wide range of depths [1]. At the later stages ($k > 1$), the adjusted depth ranges are centered at the estimated depth map from the previous stage, which is different for each pixel. We then rewrite (1) to define the sampled depth plane $d_i^k(p)$ for a pixel $p$ as follows:

$$d_i^k(p) = d_{min}^k(p) + i\Delta_k, \quad i = 1, .., M_k \tag{2}$$

$$d_{min}^k(p) = \hat{D}_q^{k-1}(p) - \frac{M_k \Delta_k}{2} \tag{3}$$

where $\hat{D}_q^{k-1}(p)$ is the predicted depth value of the pixel $p$ from the last stage. Instead of having a constant minimum depth value $d_{min}^k$, we leverage $\hat{D}_q^{k-1}$ to obtain adaptive $d_{min}^k(p)$ for each pixel $p$. The adaptive $d_{min}^k(p)$ narrows the sampled depth ranges and allows RGBD-Net to produce more accurate depth maps. The width and height of the PSVs are doubled when $k$ is increased by one (see Fig. 2). Therefore, we set $M_k = M_{k-1}/2$ and $\Delta_k = \Delta_{k-1}/2$ to avoid memory issues. For more details related to the depth plane resampling, please see the supplementary material.

**Multi-scale depth loss.** We train our depth network $R$ using the scaled depth loss function $\mathcal{L}_d = \sum_{k=1}^K \lambda_k L_k$ where $L_k$ refers to the ground-truth depth loss at the $k$-th stage and $\lambda_k$ refers to its corresponding loss weight. The depth loss $L_k$ is a combination of gradient, normal and L1 depth loss

---

$^1 d_{min}^1$, $M_1$ and $\Delta_1$ are chosen differently for each dataset

in logarithmic scale [21]. Note that, the depth loss $\mathcal{L}_d$ is optional and we show that our proposed method can be trained end-to-end using only the image loss in Section 3.3.

### 3.2. Depth-aware refinement network $G$

This section introduces a depth-aware refinement network $G$ to render the novel images based on the multi-scale depth prediction. We use differentiable bilinear interpolation $Warp$, from Jaderberg et al. [23], to map $N$ reference views $I_n$ to the target view $\hat{I}_n^k$ at the pose $s_q$, using the previously predicted depth map $\hat{D}_q^k$. We then combine the set $\Omega = \{\hat{I}_n^k\}_{n=1,...,N}$ of the warped novel images to obtain the unified warped image $\hat{W}_q^k$ as follows:

$$\hat{I}_n^k = Warp(I_n, \hat{D}_q^k, s_n, s_q) \tag{4}$$

$$\hat{W}_q^k = C(\Omega, \hat{D}_q^k) \tag{5}$$

where $C$ is the combination function which dynamically gives more weight to the reference views near the target view. Supplementary material provides more details on the combination function $C$.

Each multi-scale warped novel image $\hat{W}_q^k$ might contain several warping artifacts or missing areas due to occlusions. To address these issues, we utilize a 2D U-Net based convolutional architecture. Figure 3 shows the illustration of our proposed depth-aware refinement network $G$ using $K = 3$ scales. The high resolution warped novel view is first fed to the network to produce lower-resolution feature maps. The subsequent blocks in the encoder are reinforced by the coarser warped novel view to leverage multi-scale features. Also, multiple Dense Multi-scale Fusion Blocks (DMFB) [39] have been employed in the skip connections between the encoder and decoder to fill the missing pixels of the warped novel images. Each block adopts the combination and fusion of hierarchical features extracted from
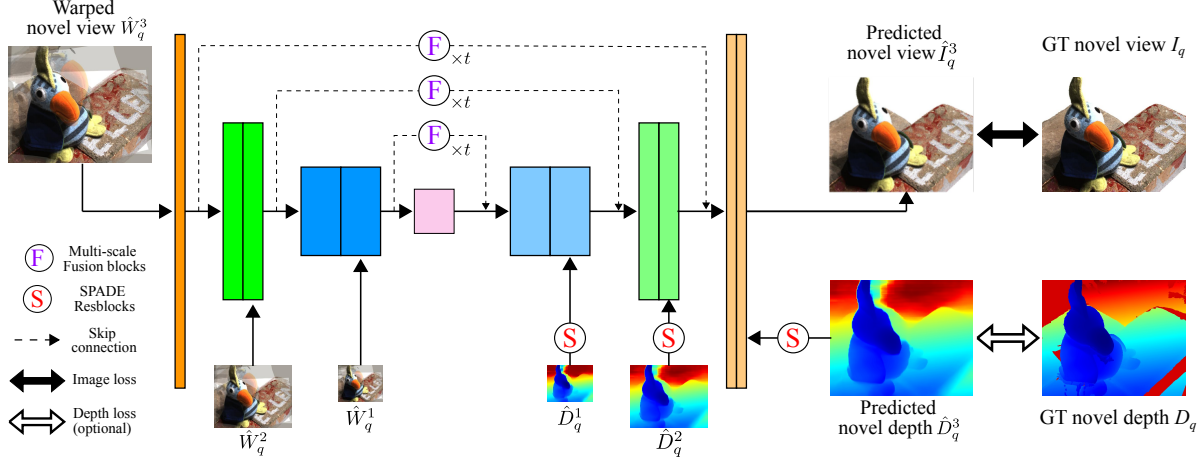
Figure 3. The architecture of the Depth-aware refinement network $G$. It refines the multi-scale warped novel views $\hat{W}_q^k$ using a U-Net architecture to produce multi-scale target images $\hat{I}_q^k$. Shown for three scales, i.e. $K = 3$ which is used in all experiments.

various convolutions with different dilation rates to obtain better multi-scale features.

To further enhance the overall quality of the predicted novel image, we leverage the complete predicted depth maps $\hat{D}_q^k$. Recent successes on conditional image synthesis [40, 61, 26] have shown that we can generate photo-realistic images conditioned on certain image data such as an image from another domain or a semantic segmentation map. We observe that our proposed depth regression network $R$ produces depth maps with sharp edges. Therefore, we exploit this to guide the image synthesizing model to produce sharp novel images. In Fig. 3, the multi-scale predicted depth maps $\hat{D}_q^k$ are fed to the Spatially-Adaptive Denormalization (SPADE) Resblocks [40] to progressively predict the novel views in a coarse-to-fine manner. More specifically, each SPADE Resblock regularizes the decoder's learned features based on the multi-scale depth predictions. This transformation ensures that the predicted novel image has similar sharp edges as the depth map produced by the $R$ network. Besides, utilizing the depth-maps in the SPADE blocks provides a valid inductive bias to the refinement network for synthesizing the novel views even without explicitly learning the depth maps with ground truth.

As can be seen in Fig. 2 and Fig. 3, the proposed method produces multi-scale depth maps and novel images at the target camera. To train both $R$ and $G$ network, the image loss between the predicted and ground-truth novel view is required. Since we are using multi-stereo datasets, the ground-truth depth maps are also available. This allows us to train our method with or without the depth loss.

### 3.3. Training

**Learning objective.** We trained the proposed method with an L1 image loss $\mathcal{L}_{l1}$, perceptual loss $\mathcal{L}_p$ [6] and hinge GAN

loss $\mathcal{L}_G$ [14] between the generated and ground-truth novel image. If the ground-truth depth map is available we can also use the scaled depth loss $\mathcal{L}_d$. The total loss is then $\mathcal{L}_{total} = \lambda_{l1}\mathcal{L}_{l1} + \lambda_p\mathcal{L}_p + \lambda_G\mathcal{L}_G + \lambda_d\mathcal{L}_d$. Note that our method does not strictly need the depth loss $\mathcal{L}_d$, which enables training on datasets that do not have ground-truth depth maps.

**Implementation details.** The models were trained with the Adam optimizer using a 0.004 learning rate for the discriminator, 0.001 for both the depth regression $R$ and refinement generator $G$ and momentum parameters (0,0.9). $\lambda_{l1} = 1, \lambda_p = 10, \lambda_{GAN} = 1, \lambda_d = 1, K = 3, N = 4, t = 5, W = 640, H = 512$. We implemented RGBD-Net in PyTorch [41], and training took 2-3 days on 4 Tesla V100 GPUs. More details are available in the supplementary material.

## 4. Experiments

We first describe our experimental setup in Section 4.1. We then evaluate RGBD-Net against the current top-performing view synthesis methods in Section 4.2. Furthermore, we additionally evaluate our predicted depth map against those produced by Multi-View Stereo (MVS) methods in Section 4.3.

### 4.1. Experimental settings

**Source image selections.** We follow the view selection method of MVSNet [56] to select the top 10 closest source images to each target image. During training, we randomly select source images among the 10 closest views as inputs to our method. Also, to reduce GPU memory requirements we use a fixed $N$ input images at each training step.

**Datasets.** We train RGBD-Net using the DTU [1] and

| Methods | Tanks & Temples [27] | | | DTU [1] | | | BlendedMVS [57] | | |
|---|---|---|---|---|---|---|---|---|---|
| | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ | LPIPS↓ | SSIM↑ | PSNR↑ |
| EVS [10] | 0.57 | 0.49 | 13.06 | 0.61 / 0.53 | 0.938 / 0.917 | 23.07 / 21.23 | 0.67 | 0.786 | 17.29 |
| LLFF [33] | 0.64 | 0.47 | 11.25 | 0.51 / 0.44 | 0.939 / 0.926 | 22.44 / 22.43 | 0.56 | 0.794 | 18.21 |
| FVS [44] | 0.18 | 0.86 | 20.26 | 0.25 / 0.30 | 0.972 / 0.951 | 26.96 / 24.08 | 0.25 | 0.815 | 22.94 |
| NPBG [2] | 0.24 | 0.82 | 19.46 | 0.36 / 0.42 | 0.942 / 0.945 | 24.78 / 23.91 | 0.36 | 0.801 | 20.18 |
| NeRF [34] | 0.62 | 0.58 | 15.69 | 0.24 / 0.29 | 0.976 / 0.958 | 28.21 / 25.81 | 0.28 | 0.824 | 22.51 |
| RGBD-Net* | 0.17 | 0.88 | 20.35 | 0.20 / 0.29 | 0.980 / 0.962 | 31.69 / 26.18 | 0.21 | 0.838 | 23.52 |
| **RGBD-Net** | **0.16** | **0.89** | **21.28** | **0.19 / 0.27** | **0.985 / 0.971** | **32.65 / 26.49** | **0.18** | **0.859** | **25.13** |

Table 1. Comparison with state-of-the-art view synthesis methods. For the DTU [1] dataset, the metrics (average over all images) are reported separately for view interpolation / extrapolation. For other datasets, we calculate the average over all target views. The RGBD-Net is trained with combined loss $\mathcal{L}_{total}$, achieving best results. Performance without the depth loss $\mathcal{L}_d$, denoted RGBD-Net*, is competitive.
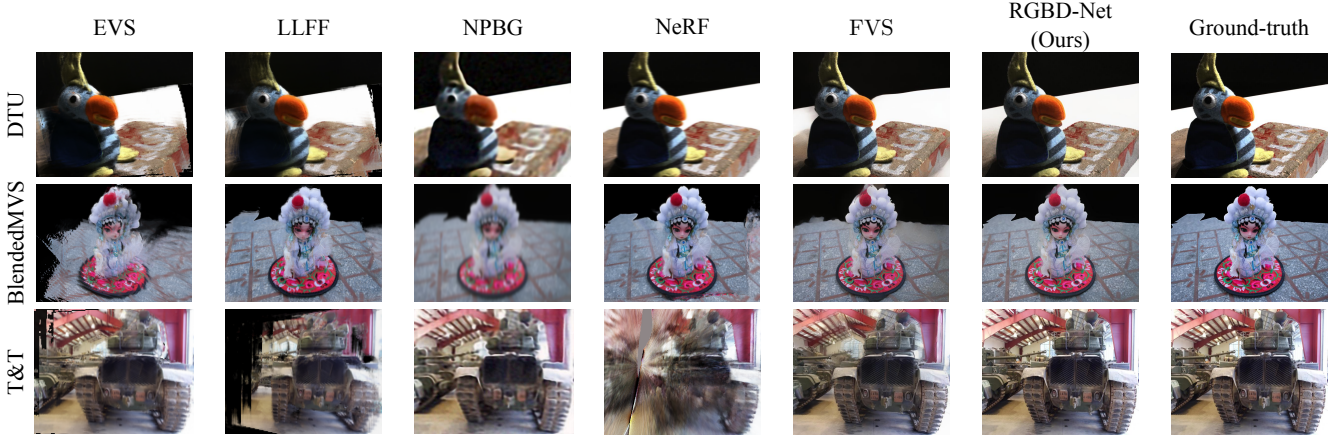


Figure 4. Novel views generated by RGBD-Net and other state-of-the-art methods in selected scenes from form the DTU [1], BlendedMVS [57] and the "intermediate" Tanks and Temples (T&T) [27] datasets.

BlendedMVS [57] datasets. DTU is an MVS dataset consisting of more than 100 scenes scanned in 7 different lighting conditions at 49 positions. From 49 camera poses, we selected 10 as targets for view synthesis and used the rest for source image selection. View extrapolation is tested using 4 target views on the extreme camera positions and we test the view interpolation on the remaining 6 views. Blended-MVS [57] is another large-scaled MVS dataset which contains high quality rendered and real images with realistic ambient lighting. We first train RGBD-Net on DTU training set and then fine-tune it on the BlendedMVS training set. We found that this yields better performance than training from scratch. We evaluate the performance of our model using the DTU and BlendedMVS testing sets.

To test the generalization capability of our method, we trained RGBD-Net on BlendedMVS and test it on the intermediate set of the Tanks and Temples [27] dataset. This set consists of 8 scenes including Family, Francis, Horse, Lighthouse, M60, Panther, Playground, and Train. More specifically, we randomly sample a subset of camera poses from each scene and treat them as the targets to be predicted.

We again note that our method has not seen these images from the Tanks and Temples dataset during training.

## 4.2. Novel view evaluation

**Baselines.** In this evaluation, we compare our approach to recently proposed methods on novel view synthesis. Among the current state-of-the-art view synthesis methods, Neural Radiance Fields (NeRF) [34] and Neural Point-based Graphics (NPBG) [2] require retraining on the test scenes, whereas our method does not need any adaptation or fine-tuning on new scenes. We use their public training code to train their models for each test scene. We also compare our method with three image-based rendering methods: LLFF [33], EVS [10] and FVS [44]. Local Light Field Fusion (LLFF) is based on the MPI representation and assumes that camera poses lie on the same plane. Extreme View Synthesis focuses on extreme stereo baseline magnification and utilizes DeepMVS [22] for depth prediction. We use the provided pretrained LLFF and EVS models to produce novel images and compare them with ours. For a fair evaluation we compare our results with the IBR meth-

ods FVS [44], LLFF [33] and EVS [10] using the same sets of unstructured source images. Lastly, in order to compare our results with NPBG [2] and NeRF [34], we perform per-scene training again using the same sets of source images.

**Metrics.** We report the PSNR, SSIM, and perceptual similarity (LPIPS) [59] of view synthesis between RGBD-Net and other state-of-the-art methods.

**Results.** We summarize the quantitative and qualitative results in Table 1 and Fig. 4. We observe that RGBD-Net outperforms the baseline methods. Also, the model RGBD-Net*, which is trained without the ground truth depth loss $\mathcal{L}_d$, shows almost similar performance to the full model, while still being better than other baselines. We also observe no significant differences between the predicted novel views produced by RGBD-Net when trained with or without depth supervision. The goal of view synthesis is to produce faithful novel views and for that purpose, we do not strictly need the predicted depth map at the target view to be perfectly accurate. More qualitative results on estimated depth maps are in the supplementary material.

We evaluate the generalization of RGBD-Net against state-of-the-art methods on the intermediate set of the Tanks and Temples dataset. The MPI-based method LLFF [33] fails to produce the complete novel image in this unstructured setting. In LLFF [33], the input views are captured within a small baseline similar to the light-field setup. Therefore, the assumptions of LLFF are not met, which leads to both missing pixels and ghosting artifacts. The depth-warp-based method EVS [10] produces slightly better results but often misses the fine details of the objects or a part of the generated image. RGBD-Net also shows better reconstruction results than FVS [44]. We found that the current state-of-the-art method on view synthesis, NeRF [34], struggles to produce good results on this dataset. The predicted novel views are either blurry or fail completely as can be seen in Fig. 4. This may be explained by the fact that NeRF makes an assumption of inward-facing scenes while our method does not have such limitation. In addition, RGBD-Net can generate plausible novel views in a wide variety of 3D scenes. Comparing to NPBG [2], our method produces sharper novel views and shows superior performance in terms of the LPIPS metrics as can be seen in the Table 1.

We then evaluate our method on the multi-view stereo DTU [1] and BlendedMVS [57] test sets. Different from the Tanks and Temples dataset, these multi-view stereo datasets contain in-ward facing images. Moreover, the source camera poses are densely sampled around the target pose. Therefore, these datasets fit to the assumptions of NeRF [34] and LLFF [33]. As can be seen in Table 1 and Fig 4, the results of LLFF and NeRF are clearly better compared to their performances on the Tanks and Temples dataset. However, we still notice blending artifacts in the
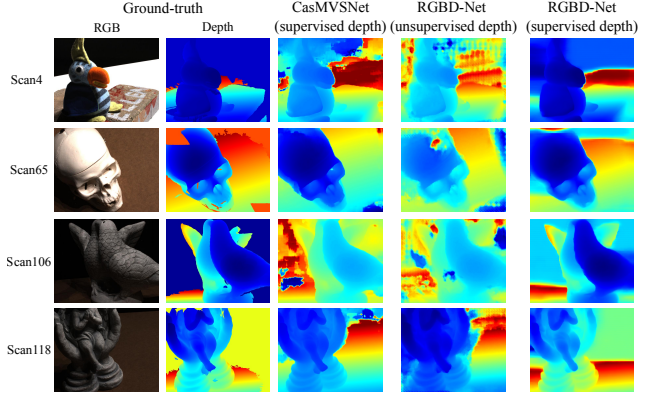


Figure 5. Depth map estimates by RGBD-Net and the currently top performing CasMVSNet [17] on selected scenes from DTU[1].

novel views produced by EVS [10], which reflect the lower quantitative performance. NPBG [2] produces blurry target views that often miss fine details. Despite of the impressive synthesized foreground in the target views, FVS produces blurry results on the background. In contrast, RGBD-Net leverages the estimated complete multi-scale depth maps to guide the generator network $G$ and synthesize also plausible background.

## 4.3. Multi-view stereo evaluation

**Baselines.** We evaluate the predicted depth maps produced by RGBD-Net against MVS methods [56, 7, 58, 9, 17] on the DTU [1] test set and the *intermediate* set of Tanks and Temples [27] dataset. We use fusible [13] as the post-processing step to reconstruct the 3D point cloud of the scene. Therefore, a more accurate and consistently estimated depth map would lead to better performance in 3D reconstruction. We note that our problem is a more challenging case since RGBD-Net predicts both the target depth maps and color images using only the reference views. Whereas, all other MVS methods predict only the depth maps at the reference poses while using the reference images as input.

**Metrics.** We calculate the mean *accuracy*, *completeness* and *overall* using the evaluation code provided by [1]. The average of mean accuracy and completeness represent the reconstruction quality. Moreover, we also calculate the mean F-score on the Tanks and Temples [27] dataset.

**Results.** As can be seen in Table 2, the proposed method trained without the depth loss (RGBD-Net*) shows competitive performance with other MVS baselines. Using the ground-truth depth loss, our full model (RGBD-Net) achieves the best mean F-score on Tanks and Temples and the best overall distance on DTU. In Fig. 5, we show qualitative results on the predicted depth map of the reference camera compared to those produced by the MVS baselines.

| | DTU [1] | | | Tanks and Temples [27] | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Acc.↓ (mm) | Comp.↓ (mm) | Overall↓ (mm) | Mean↑ | Family↑ | Francis↑ | Horse↑ | Lighthouse↑ | M60↑ | Panther↑ | Playground↑ | Train↑ |
| MVSNet [56] | 0.456 | 0.646 | 0.551 | 43.48 | 55.99 | 28.55 | 25.07 | 50.79 | 53.96 | 50.86 | 47.90 | 34.69 |
| Point-MVSNet [7] | 0.361 | 0.421 | 0.391 | 48.27 | 61.79 | 41.15 | 34.20 | 50.79 | 51.97 | 50.85 | 52.38 | 43.06 |
| PVA-MVSNet [58] | 0.352 | 0.414 | 0.383 | 54.46 | 69.36 | 46.80 | 46.01 | 55.74 | 57.23 | 54.75 | 56.70 | 49.06 |
| UCSNet [9] | 0.330 | 0.392 | 0.361 | 54.83 | 76.09 | 53.16 | 43.04 | 54.00 | 55.60 | 51.49 | 57.38 | 47.89 |
| CasMVSNet [17] | 0.325 | 0.385 | 0.355 | 56.84 | 76.37 | 58.45 | 46.26 | 55.81 | 56.11 | 54.06 | 58.18 | 49.51 |
| RGBD-Net* | 0.334 | 0.390 | 0.362 | 55.57 | 76.82 | 53.84 | 44.29 | 55.02 | 55.98 | 52.78 | 58.63 | 47.25 |
| **RGBD-Net** | **0.320** | **0.381** | **0.349** | **59.32** | **77.01** | **60.25** | **47.09** | **63.45** | **62.19** | **55.16** | **59.27** | **50.19** |

Table 2. Point cloud accuracy on the DTU test [1] and Tanks and Temples [27] *intermediate* datasets.

We observe that our method is able to generate accurate depth map of the target view without using the reference image as input. Moreover, we also show that our method performs well on unseen data. As can be seen in the Table 2 and Fig. 1, our proposed method is able to predict both the depth maps and color images of target views and then use them to reconstruct the 3D point cloud on the Tanks and Temples [27] dataset. More examples of the generated point cloud are in the the supplementary material.
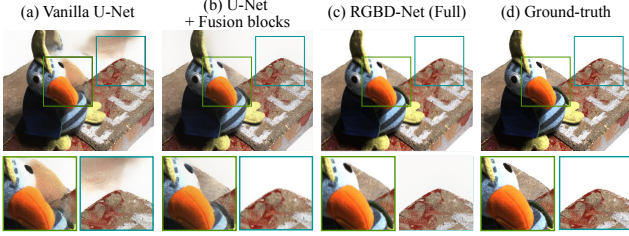


Figure 6. Comparison of the ground-truth novel image (d) with predicted novel views by: (a) Vanilla U-Net model, (b) U-Net with multi-scale fusion blocks (DMFB) and (c) RGBD-Net full model.

| U-Net | Depth scales | DMFB | SPADE | LPIPS↓ | SSIM↑ | PSNR↑ |
|---|---|---|---|---|---|---|
| ✓ | | | | 0.42 | 0.927 | 23.89 |
| ✓ | ✓ | | | 0.31 | 0.941 | 27.44 |
| ✓ | ✓ | ✓ | | 0.26 | 0.967 | 28.67 |
| ✓ | ✓ | ✓ | ✓ | **0.235** | **0.972** | **29.21** |

Table 3. RGBD-Net architecture ablation study. Reconstruction accuracy of novel view synthesis on the DTU test set [1].

| | # of reference images | | | | | | |
|---|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
| LPIPS↓ | 0.450 | 0.362 | 0.261 | **0.235** | 0.235 | 0.247 | 0.252 |
| SSIM↑ | 0.935 | 0.949 | 0.969 | **0.975** | 0.975 | 0.965 | 0.961 |

Table 4. The impact of the number of reference image. Measured by reconstruction accuracy of novel view synthesis on the DTU test set [1].

Finally, we report the average execution time of our model tested on an RTX 2080 GPU. To render a color and depth image with the size of $640 \times 512$ using 4 reference views, the depth regression network $R$ takes 102 ms and the refinement network $G$ takes 54 ms, respectively. Comparisons on the average execution time of RGBD-Net and other methods are included in the supplementary material.

## 5. Ablation study

In this section, we show how each proposed module affects the overall performance of RGBD-Net. Table 3 summarizes the quantitative results on various architecture choices using the test set of the DTU dataset [1]. We notice that training the single scale depth regression network and the vanilla U-Net refinement network does not produce plausible target views as they contain cluttered background and incorrect geometry. To address this issue, we utilize multiple DMFB which allows RGBD-Net to have bigger receptive field. As can be seen in Fig. 6, the overall quality of the predicted novel views improves significantly by combining multi-scale depth predictions and DMFB. However, the predicted novel views are often blurry compared to the ground-truth images. To further enhance the quality of the generated images, our full model employs both multiple DMFB and SPADE Resblocks in the generator $G$. Table 3 shows that SPADE Resblocks have a substantial impact on the results and they clearly help our method to achieve state-of-the-art performance.

In Table 4, we evaluate the performance of our method with an increasing number of source images using the DTU [1] dataset. We report both SSIM and LPIPS metrics with the number of source images up to 7. We observe that RGBD-Net performs the best with 4 input views and then the results start saturating. More reference views are required if the target pose is farther away or there is less overlap with the viewing frustum of the target image.

## 6. Conclusions

We presented RGBD-Net, a method for novel view synthesis in the challenging setting of unstructured input views. Our method first predicts the novel depth maps in a coarse-to-fine manner and then refines the depth-warped novel views to produce the final target images. Moreover, RGBD-Net outperformed state-of-the-art view synthesis models

and it also showed competitive performance on the 3D reconstruction task compared to recent MVS methods.

# References

[1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *Int. J. Comput. Vision*, 120(2):153–168, Nov. 2016. 2, 5, 6, 7, 8, 12, 13, 15, 18

[2] Kara-Ali Aliev, Dmitry Ulyanov, and Victor S. Lempitsky. Neural point-based graphics. *CoRR*, abs/1906.08240, 2019. 2, 3, 6, 7, 12

[3] Michael Broxton, Jay Busch, Jason Dourgarian, Matthew DuVall, Daniel Erickson, Dan Evangelakos, John Flynn, Peter Hedman, Ryan Overbeck, Matt Whalen, and Paul Debevec. Deepview immersive light field video. In *ACM SIGGRAPH 2020 Immersive Pavilion*, SIGGRAPH '20, New York, NY, USA, 2020. Association for Computing Machinery. 2

[4] 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*, SIGGRAPH '01, page 425–432, New York, NY, USA, 2001. Association for Computing Machinery. 1, 2

[5] Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM Trans. Graph.*, 32(3), July 2013. 1, 2

[6] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 1520–1529. IEEE Computer Society, 2017. 5

[7] Rui Chen, Songfang Han, Jing Xu, and Hao Su. Point-based multi-view stereo network. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1538–1547, 2019. 7, 8

[8] Shenchang Eric Chen and Lance Williams. View interpolation for image synthesis. In *Proceedings of the 20th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '93, page 279–288, New York, NY, USA, 1993. Association for Computing Machinery. 1, 2

[9] Shuo Cheng, Zexiang Xu, Shilin Zhu, Zhuwen Li, Li Erran Li, Ravi Ramamoorthi, and Hao Su. Deep stereo using adaptive thin volume representation with uncertainty awareness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2524–2534, 2020. 7, 8

[10] Inchang Choi, Orazio Gallo, Alejandro Troccoli, Min H Kim, and Jan Kautz. Extreme view synthesis. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7781–7790, 2019. 1, 3, 6, 7, 12

[11] R. 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, 1996. 2

[12] John Flynn, Michael Broxton, Paul Debevec, Matthew DuVall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2367–2376, 2019. 2, 3

[13] S. Galliani, K. Lasinger, and K. Schindler. Massively parallel multiview stereopsis by surface normal diffusion. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 873–881, 2015. 7, 13

[14] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems 27*, pages 2672–2680. Curran Associates, Inc., 2014. 5

[15] 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*, SIGGRAPH '96, page 43–54, New York, NY, USA, 1996. Association for Computing Machinery. 2

[16] N. Greene. Environment mapping and other applications of world projections. *IEEE Computer Graphics and Applications*, 6(11):21–29, 1986. 2

[17] 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. 3, 7, 8, 12, 13

[18] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM Trans. Graph.*, 37(6), Dec. 2018. 2

[19] Peter Hedman, Tobias Ritschel, George Drettakis, and Gabriel Brostow. Scalable inside-out image-based rendering. *ACM Trans. Graph.*, 35(6), Nov. 2016. 2

[20] Benno Heigl, Reinhard Koch, Marc Pollefeys, Joachim Denzler, and Luc J. Van Gool. Plenoptic modeling and rendering from image sequences taken by hand-held camera. In *Mustererkennung 1999, 21. DAGM-Symposium*, page 94–101, Berlin, Heidelberg, 1999. Springer-Verlag. 2

[21] Junjie Hu, Mete Ozay, Yan Zhang, and Takayuki Okatani. Revisiting single image depth estimation: Toward higher resolution maps with accurate object boundaries. In *IEEE Winter Conference on Applications of Computer Vision, WACV 2019, Waikoloa Village, HI, USA, January 7-11, 2019*, pages 1043–1051. IEEE, 2019. 4

[22] Po-Han Huang, Kevin Matzen, Johannes Kopf, Narendra Ahuja, and Jia-Bin Huang. Deepmvs: Learning multi-view stereopsis. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 3, 6

[23] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'15, page 2017–2025, Cambridge, MA, USA, 2015. MIT Press. 4

[24] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM Trans. Graph.*, 35(6), Nov. 2016. 2

[25] Animesh Karnewar and Oliver Wang. Msg-gan: Multi-scale gradients for generative adversarial networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7799–7808, 2020. 3

[26] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020. 5

[27] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 2, 6, 7, 8, 12, 13, 14

[28] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '96, page 31–42, New York, NY, USA, 1996. Association for Computing Machinery. 1, 2

[29] Zhengqi Li, Wenqi Xian, Abe Davis, and Noah Snavely. Crowdsampling the plenoptic function. In *Proc. European Conference on Computer Vision (ECCV)*, 2020. 3

[30] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2117–2125, 2017. 3

[31] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *NeurIPS*, 2020. 1, 3

[32] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4), July 2019. 1, 3

[33] Ben Mildenhall, Pratul P. Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 2019. 1, 3, 6, 7, 12

[34] 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. In *ECCV*, 2020. 1, 2, 3, 6, 7, 12, 13

[35] Hans-Peter Seidel Mojtaba Bemana, Karol Myszkowski and Tobias Ritschel. X-fields: Implicit neural view-, light- and time-image interpolation. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia 2020)*, 39(6), 2020. 2, 3

[36] Thu Nguyen-Phuoc, Chuan Li, Lucas Theis, Christian Richardt, and Yong-Liang Yang. Hologan: Unsupervised learning of 3d representations from natural images. In *The IEEE International Conference on Computer Vision (ICCV)*, Nov 2019. 3

[37] Thu Nguyen-Phuoc, Christian Richardt, Long Mai, Yong-Liang Yang, and Niloy Mitra. Blockgan: Learning 3d object-aware scene representations from unlabelled images. *arXiv preprint arXiv:2002.08988*, 2020. 3

[38] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3

[39] Evangelos Ntavelis, Andrés Romero, Siavash Bigdeli, Radu Timofte, et al. AIM 2020 challenge on image extreme inpainting. In *European Conference on Computer Vision Workshops*, 2020. 4

[40] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 5

[41] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8026–8037. Curran Associates, Inc., 2019. 5

[42] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Trans. Graph.*, 36(6), Nov. 2017. 1, 2

[43] Thomas Porter and Tom Duff. Compositing digital images. In *Proceedings of the 11th Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '84, page 253–259, New York, NY, USA, 1984. Association for Computing Machinery. 3

[44] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *European Conference on Computer Vision*, 2020. 3, 6, 7, 12

[45] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 3

[46] Steven M. Seitz and Charles R. Dyer. View morphing. In *Proceedings of the 23rd Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '96, page 21–30, New York, NY, USA, 1996. Association for Computing Machinery. 2

[47] N. K. Shukia, S. Sengupta, and M. Chakraborty. Intermediate view synthesis in wide-baseline stereoscopic video for immersive telepresence. In *The 2nd European Workshop on the Integration of Knowledge, Semantics and Digital Media Technology, 2005. EWIMT 2005. (Ref. No. 2005/11099)*, pages 83–88, 2005. 1

[48] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2437–2446, 2019. 2

[49] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2019. 3

[50] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in*

*Neural Information Processing Systems*, pages 1121–1132, 2019. 1, 2

[51] P. P. Srinivasan, R. Tucker, J. T. Barron, R. Ramamoorthi, R. Ng, and N. Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 175–184, 2019. 3

[52] A. Tewari, O. Fried, J. Thies, V. Sitzmann, S. Lombardi, K. Sunkavalli, R. Martin-Brualla, T. Simon, J. Saragih, M. Nießner, R. Pandey, S. Fanello, G. Wetzstein, J.-Y. Zhu, C. Theobalt, M. Agrawala, E. Shechtman, D. B Goldman, and M. Zollhöfer. State of the Art on Neural Rendering. *Computer Graphics Forum (EG STAR 2020)*, 2020. 2

[53] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2

[54] Justus Thies, Michael Zollhöfer, Christian Theobalt, Marc Stamminger, and Matthias Nießner. Image-guided neural object rendering. In *International Conference on Learning Representations*, 2019. 2

[55] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 3

[56] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *European Conference on Computer Vision (ECCV)*, 2018. 3, 5, 7, 8, 12, 13

[57] Yao Yao, Zixin Luo, Shiwei Li, Jingyang Zhang, Yufan Ren, Lei Zhou, Tian Fang, and Long Quan. Blendedmvs: A large-scale dataset for generalized multi-view stereo networks. *Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 6, 7, 12, 13, 16, 17, 19

[58] Hongwei Yi, Zizhuang Wei, Mingyu Ding, Runze Zhang, Yisong Chen, Guoping Wang, and Yu-Wing Tai. Pyramid multi-view stereo net with self-adaptive view aggregation. *arXiv preprint arXiv:1912.03001*, 2019. 7, 8

[59] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 7

[60] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. In *SIGGRAPH*, 2018. 1, 3

[61] Peihao Zhu, Rameen Abdal, Yipeng Qin, and Peter Wonka. Sean: Image synthesis with semantic region-adaptive normalization. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 5

# Appendix of RGBD-Net

The organization of the appendix is as follows. Section A shows the comparisons on the average execution time of our method and other view synthesis methods. We present the function $C$ to combine multiple warped novel views in Section B. More implementation details about RGBD-Net are provided in Section C. Finally, additional qualitative results are presented in Section D.

## A. Execution time

| Methods | EVS [10] | LLFF [33] | FVS [44] | NPBG [2] | NeRF [34] | RGBD-Net (ours) |
|---|---|---|---|---|---|---|
| Avg time (ms / img) | 984 | 651 | 541 | 352 | 8153 | **156** |

Table 5. Comparisons on the average execution time of RGBD-Net and other view synthesis methods.

In Table. 5, we report the average execution times to synthesize a novel image between RGBD-Net with different methods. In all experiments, we synthesize the novel image with the size of $640 \times 512$ pixels using 4 reference views on the T&T datasets. Notice that, RGBD-Net is 52 times faster than the current state-of-the-art neural rendering method NeRF [34]. Moreover, our method also run faster than other image-based rendering techniques while maintaining superior performance.

## B. Combining warped novel views

In section 3.2 of the paper, we present a combination function $C$, which dynamically gives more weight to the reference views near the target view. $C$ in equation (5) of the paper is a merging function from the set $\Omega = \{\hat{I}_n^k\}_{n=1,...,N}$ of the warped novel images to obtain the unified warped image $\hat{W}_q^k$ as follows:

$$\hat{W}_q^k = (\sum_{n=1}^{N} \omega_{d_n}\omega_{a_n})^{-1} \sum_{n=1}^{N} \omega_{d_n}\omega_{a_n}\hat{I}_n^k \qquad (6)$$

where $\omega_{d_n}$ and $\omega_{a_n}$ are the distance and angle weight of the $n$-th reference pose with respect to the target camera, respectively. We define $\omega_{d_n}$ and $\omega_{a_n}$ as follows:

$$\omega_{d_n} = \exp(\frac{-\|p_n - p_q\|_1}{\sigma_d}) \qquad (7)$$

$$\omega_{a_n} = \exp(\frac{-|v_n v_q^{\intercal}|_1}{\sigma_v}) \qquad (8)$$

where $p_n$ and $p_q$ are the translation vectors, $v_n$ and $v_q$ are the rotation matrices of the $n$-th reference and target poses, respectively. We set $\sigma_d = 0.04$ and $\sigma_v = 0.16$ as scaling factors.

## C. Implementation details

In this section, we first provide implementation details how to train RGBD-Net. We then explain how we define hypothesis depth planes using our proposed depth plane resampling. Next we show how the depth regression network $R$ produces multi-scale depth map of the target view. Finally, we provide more details on how to generate point-cloud using RGBD-Net.

**Training.** We trained RGBD-Net on the DTU dataset for 25 epochs and subsequently finetuned it on the Blended-MVS dataset for 30 epochs. The models were trained with the Adam optimizer with a batch size of 4. In both training sets, the input and output reference views have the same image size of $640 \times 512$ pixels and we set the number of reference views to $N = 4$. To balance between accuracy and efficiency, we adopt a three-scale ($K = 3$) depth regression network $R$. Accordingly, the spatial resolution of extracted feature maps $F_n^k$ is set to 1/16, 1/4 and 1 of the original image size.

**Depth plane resampling.** In Section 3.1 of the paper, the depth plane resampling depends on the initial number of hypothesis depth plane $M_1$, minimum depth value $d_{min}^1$ and depth interval $\Delta_1$. In all experiments, we use the same the number of hypothesis depth plane $M_1 = 48$ at the first stage. Both $d_{min}^1$ and $\Delta_1$ are chosen differently for each dataset.

The DTU [1] dataset is captured in a controlled environment so we follow previous works [56, 17] on MVS and set $d_{min}^1 = 425$ and $\Delta_1 = 10.6$. The BlendedMVS [57] and Tanks and Temples (T&T) [27] contain large-scaled scenes which have variety of depth ranges. Some depth ranges are from 0.1 to 2 or from 10 to 100. Note that these numbers are not the absolute distances in some known units. Therefore, we scale those depth ranges roughly to the same scale from 100 to 1000. In that way, $d_{min}^1$ and $\Delta_1$ are chosen differently per-scene based on the scaled depth range. Moreover, the scaled depth ranges in BlendedMVS and T&T datasets become approximately similar to the depth range of the DTU dataset. This allows us to train RGBD-Net on both DTU and BlendedMVS datasets and test the generalization of the model on the T&T dataset.

**Depth map regression.** Inspired by current learning-based MVS methods [56, 17], the predicted depth map is regressed from the probability volume via the *soft-argmax* operation. We denote the probability volume over all the $M_k$ depth hypothesis as $V^k$. The predicted depth value $\hat{D}_q^k(p)$ of each pixel $p$ is defined as follows:

$$\hat{D}_q^k(p) = \sum_{i=1}^{M_k} d_i^k V_i^k(p) \qquad (9)$$

where $V_i^k(p)$ is the predicted probability of the depth plane $d_i^k$ for the pixel $p$ at the $k$-th scale.
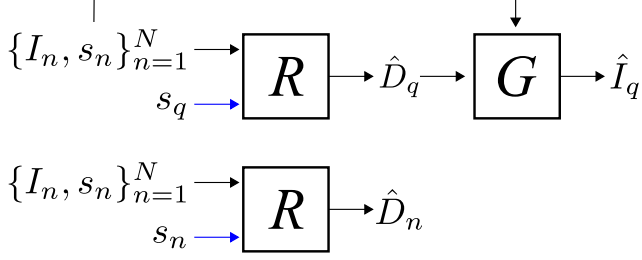
Figure 7. The illustration of how RGBD-Net produces the RGB image $\hat{I}_q$, depth map $\hat{D}_q$ at the target view $s_q$ and the depth map $\hat{D}_n$ of the reference pose $s_n$ in the testing time. Blue arrow indicates the input pose of the depth regression network $R$.

**Pointcloud generation.** Similar to previous works [56, 17] on MVS, we apply depth map filter and fusion approach [13] to merge all predicted novel views and depth maps into a unified pointcloud output. In the testing time, we predict the depth map $\hat{D}_n$ of each reference view $I_n$ using the trained depth regression network $R$ as can be seen in Fig. 7. We use the union set $\Psi = \{\hat{I}_q, \hat{D}_q\}_{q=1}^{Q} \cup \{I_n, \hat{D}_n\}_{n=1}^{N}$ of the estimated novel views and depth maps as inputs for pointcloud generation[2]. We apply two-step depth map filtering strategy to remove outliers. Finally, median depth map fusion is applied to refine all depth maps. The 3D pointcloud is obtained by projecting all refined depth maps into the 3D space. We provide qualitative results of generated pointclouds in the Section D.2.

## D. Additional qualitative results

### D.1. Novel view synthesis

In this section, we provide additional qualitative results. Fig. 8, 9 and 10 show more examples of rendered novel views using RGBD-Net and other view synthesis methods on the Tanks and Temples [27], DTU [1] and BlendedMVS [57] datasets, respectively. For NeRF [34], we manually define the bounding volume around the main object in each testing scene.

### D.2. Reconstructed pointcloud

Fig. 11, 12 and 13 show more examples of generated pointclouds using the proposed RGBD-Net on the Tanks and Temples [27], DTU [1] and BlendedMVS [57] datasets,

---

[2]We use the predictions of the final scale so superscript $k$ is omitted.

Figure 8. Additional qualitative results on Tanks and Temples dataset [27]. RGBD-Net* is our proposed RGBD-Net trained without the ground-truth depth loss. We observer no significant difference on the performance of view synthesis between RGBD-Net and RGBD-Net*.

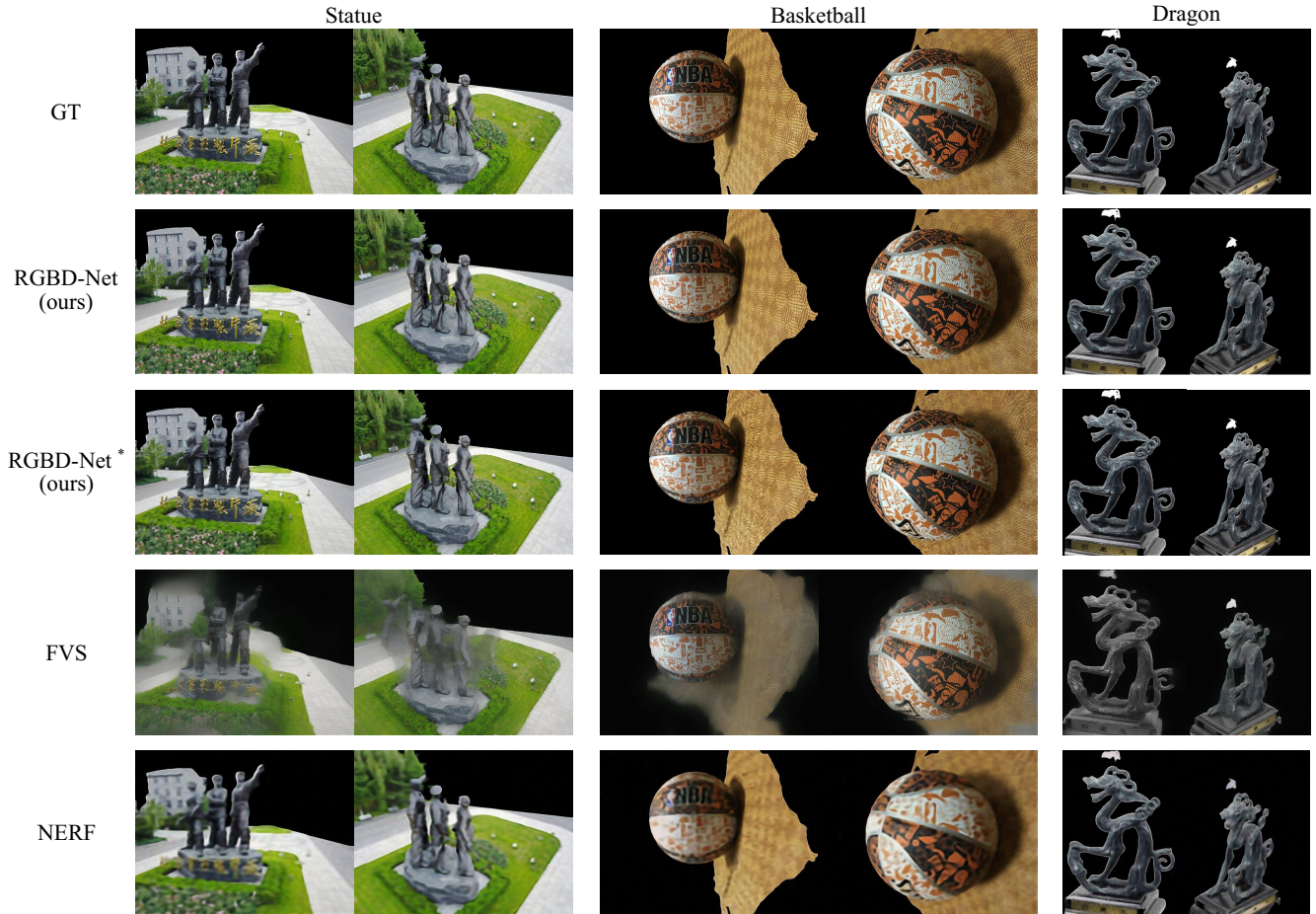Figure 9. Additional qualitative results on DTU dataset [1].

Figure 10. Additional qualitative results on BlendedMVS dataset [57].

Figure 11. Pointcloud results of RGBD-Net on the *intermediate* set of Tanks and Temples [57] dataset.

Figure 12. Pointcloud results of RGBD-Net on the DTU test set [1].

Figure 13. Pointcloud results of RGBD-Net on the BlendedMVS test set [57].