# NeuralMVS: Bridging Multi-View Stereo and Novel View Synthesis

**Radu Alexandru Rosu, Sven Behnke**
Autonomous Intelligent Systems
University of Bonn, Germany
{rosu, behnke}@ais.uni-bonn.de

## Abstract

Multi-View Stereo (MVS) is a core task in 3D computer vision. With the surge of novel deep learning methods, learned MVS has surpassed the accuracy of classical approaches, but still relies on building a memory intensive dense cost volume. Novel View Synthesis (NVS) is a parallel line of research and has recently seen an increase in popularity with Neural Radiance Field (NeRF) models, which optimize a per scene radiance field. However, NeRF methods do not generalize to novel scenes and are slow to train and test. We propose to bridge the gap between these two methodologies with a novel network that can recover 3D scene geometry as a distance function, together with high-resolution color images. Our method uses only a sparse set of images as input and can generalize well to novel scenes. Additionally, we propose a coarse-to-fine sphere tracing approach in order to significantly increase speed. We show on various datasets that our method reaches comparable accuracy to per-scene optimized methods while being able to generalize and running significantly faster.

## 1  Introduction

Multi-view Stereo (MVS) recovers depth and geometry from multiple images with known camera poses. This is usually done with classical methods like COLMAP [9], Gipuma [4], or MVE [5] by searching for correspondences along epipolar lines. These algorithms lack learned components and cannot cope with challenging conditions like imperfect calibration, blurry or incomplete images, and heavy occlusion.

Recent Novel View Synthesis (NVS) methods like NeRF [7] recover geometry as a byproduct of view synthesis. Geometry is represented as a radiance field which is multi-view consistent between all images. This has the advantage of recovering more complete depth while being more robust to imperfect images than classical methods.



Figure 1: NeuralMVS processes multiple input views to synthesize a novel colored view with corresponding depth and gives an estimate for the output confidence.

Main disadvantages of NVS methods are the large processing time and a lack of generalization. They require per-scene training which can take up to several days.

Additionally, inference speed is also limited—often requiring multiple minutes to synthesize a full image together with the corresponding depth.

More recent NVS approaches like pixelNeRF [14] and IBRNet [12] improve the reconstruction speed by having a training phase which allows the model to generalize to novel scenes and thus require only little per-scene fine-tuning. Inference is still slow as the network needs to query the radiance field multiple times during synthesis and requires a very dense sampling of the view frustum in the ray-marching step.

In our approach, we leverage ideas from MVS and NVS and combine them in a new learning-based method that generalizes well to novel scenes and reaches comparable reconstructions to per-scene optimized methods while requiring only a fraction of the time.

First, our method accelerates the ray marching step by differentiable sphere tracing. While ray marching requires hundreds of samples per ray in order to achieve good accuracy, sphere tracing can reach the object surface in as little as five iterations by predicting for each ray sample its jump towards the next one.

Second, instead of tracing one ray per pixel, our approach starts by tracing on a coarse image which is iteratively refined until we reach the full resolution. This coarse-to-fine method further alleviates the labor-intensive step of finding the 3D surface of the object.

Third, we propose a new scheme for the selection of conditioning views based on a Delaunay triangulation of the input views. We show that this is temporally more stable, improves generalization and provides better results than methods based on proximity of viewing direction.

Fourth, we introduce a loss that encourages the network to output a confidence map for the novel RGBD view. These confidence values align well with parts of the image that are undersampled or occluded and may be used to inform further reconstruction or refinement methods.

In summary, our contributions are:

- a new learning-based novel view synthesis method which generalizes to unseen scenes,
- an efficient coarse-to-fine approach based on differentiable sphere tracing to recover depth with few samples conditioned on a set of input views, and
- a loss that encourages the network to output a confidence map for each novel view produced.

The general pipeline of our method is shown in Fig. 1 and the core components are detailed in Fig. 2.

## 2  Related Work

Several methods for learned MVS have gained popularity lately. MVSNet [13] proposes an end-to-end differentiable model to learn depth inference from unstructured stereo. Features are extracted from images and a dense cost volume is built using samples at regular intervals. The volume is regularized using 3D convolutions and then used to regress a depth map. In contrast in our approach, we do not define the depth samples a priori, but rather let the network learn where to sample using a differentiable sphere tracer.

The work of Darmon *et al* [3] further builds on MVSNet and shows that depth can be recovered by using a color reconstruction loss. Similarly, we only employ an RGB loss, and do not supervise the depth map as in many settings an accurate ground-truth may not be available.

Recently, NeRF [7] has gained popularity for synthesizing highly-detailed novel views. NeRFs represent the scene as a radiance field and optimize it using differentiable ray marching. The ray marching step is computationally expensive since it samples the 3D space densely at regular intervals. One main contribution of NeRF is the introduction of the positional encoding in the context of NVS that enables the model to learn high-frequency details. In order to recover the 3D surface precisely, Mildenhall et al. propose a hierarchical sampling strategy that optimizes two NeRF models: one for coarse samples and one for fine samples closer to the surface. In our work, we leverage the
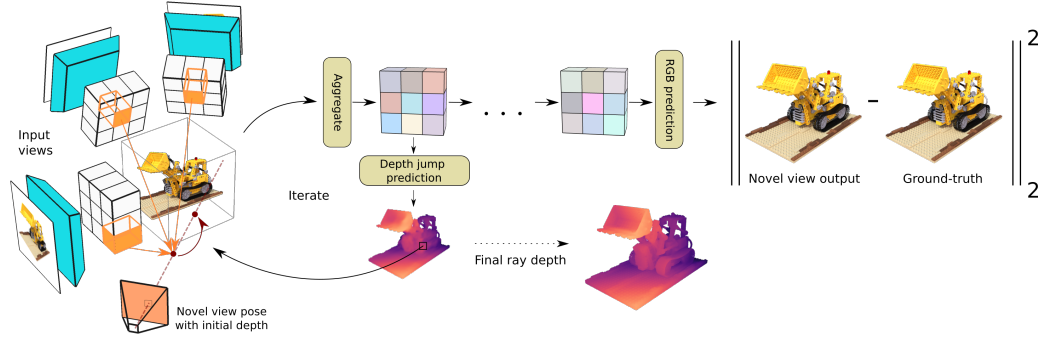
Figure 2: High-level features from the input view images are aggregated onto each ray sample. For each ray, a recursive network predicts the jump towards the next sample (red arrow). This process is repeated for a fixed number of iterations. At the last ray iteration, the final aggregated features are passed through a rendering network in order to predict the novel view RGB map. The network is only supervised with an RGB loss.

positional encoding for our ray marching step and propose sphere tracing as a method to alleviate the regular sampling of NeRF.

Scene Representation Networks (SRN) [10] is another approach which uses sphere tracing for traversing the rays in 3D space. However, this method is only able to recover low-frequency detail of the scene while we recover fine details by directly conditioning our model on the image features.

Other methods like FaDIV-Syn [8] propose to recover novel views of the scene without reconstructing depth by warping input views into the target frame at a series of predefined depth planes and letting the network learn how to best render the novel view. In contrast, we infer both the depth and the novel view jointly by explicitly letting the network modify the depth planes which are used to project input views.

## 3 Method

Given a set of source views, our method synthesizes depth and color image for a novel target view pose. The core idea is to recover a depth map for the target view such that warping source views onto it results in an image that matches the target as close as possible. Our method can thus be viewed as two jointly trained networks where the first one recovers the geometry of the scene. The second network uses that geometry as a proxy onto which the source views are projected to recover the novel view.

### 3.1 View Selection and Feature Extraction

In order to select the best suited source views to create the target view, various schemes have been proposed. Most of them are based on spatial proximity and view direction [5, 9]. However, we observe that these methods tend to fail choosing the most informative views when images are taken with non-uniform spacing. As shown in Fig. 3, we may not be able to reconstruct the whole novel view for the target position (yellow dot) depending on scene geometry when choosing the right three views (blue dots) in case of large occlusions.

In contrast, we propose to choose the "working set" based on the Delaunay triangulation of the view positions. This ensures both a better coverage of the nearby view space and allows for an easy way to compute weightings for the views by using barycentric coordinates. Since we construct the triangulation in 2D while the camera positions are in 3D, we first need to determine which view configuration is present in the scene. We distinguish between two types: hemisphere sampling in which the views are placed in the upper hemisphere around the scene and fronto-parallel sampling where they are mostly planar in front of the scene.
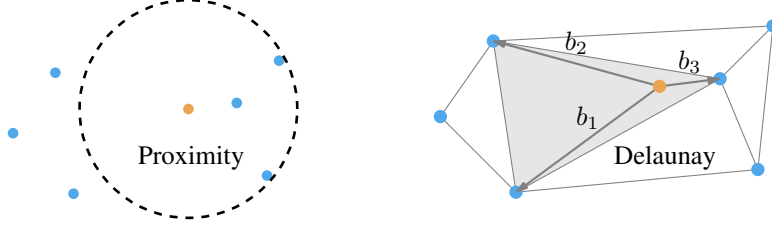
3

Figure 3: Given a novel view (orange) we need to choose a working set from all the views in the dataset (blue). Proximity-based view selection can lead to significant occlusion as the working set only views the scene from the right. Our Delaunay-based approach selects the views that belong to the closest triangle, ensuring scene viewing from different sides.
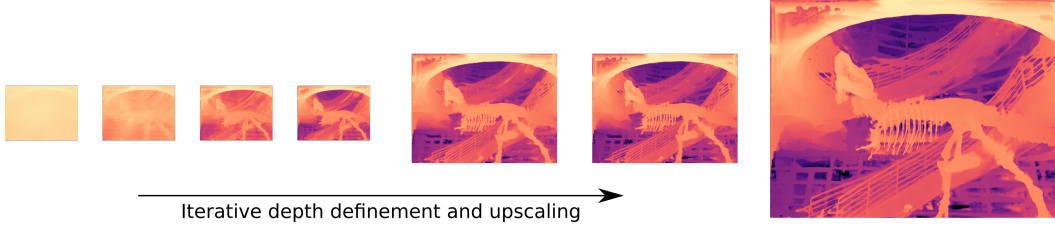


Iterative depth definement and upscaling

Figure 4: Depth estimation for the novel view starts on a coarse scale and is initialized to a constant value close to the camera. The network iteratively refines the current estimate using input view features and upscales the depth until the full resolution is reached.

For the case of hemisphere sampling, we stereographically project the camera positions onto the 2D plane where we perform the triangulation and then lift the result back to 3D. For fronto-parallel sampling, we orthographically project the camera positions onto the common plane defined by all views.

After triangulation, the closest triangle to the target view position is selected and the corresponding images form the working set. These three images are also assigned a weight $b_i$ which corresponds to the barycentric coordinate of the target view w.r.t. the triangle. Fig. 3 (right) shows selected view positions with our approach on the previous example.

Finally, we use a shared U-Net model to extract feature maps $\mathbf{F}_i \in \mathbb{R}^{H_i \times W_i \times d}$ for each image $\mathbf{I}_i$ in the working set .

## 3.2 Geometry

To recover the scene geometry, we shoot rays from each pixel of the target view and try to find the intersection of the ray with the scene surface. NeRF-like models accomplish this by densely sampling the ray at predefined intervals. Hence, most samples lay in empty space, slowing down processing. In contrast, we draw inspiration from Scene Representation Networks (SRN) [10] and propose to use a differentiable sphere tracer which predicts for each sample on the ray a jump towards the next sample. This effectively enables the network to learn and adapt the step-size, which greatly improves the sample efficiency and speed.

We parametrize each ray from the target view as follows:

$$r(t) = \mathbf{o} + t\mathbf{d}, \tag{1}$$

where $t$ is the distance along the ray, $\mathbf{o}$ is the origin of the camera in world coordinates, and $\mathbf{d}$ is the normalized direction of the ray. We initialize $t$ to be a small value such that ray marching starts close to the camera. At each ray marching step, the position of the sample $\mathbf{x} = r(t)$ is obtained in world coordinates. The raysample is projected into each source view $\mathbf{I}_i$ from where local features $\mathbf{f}_i \in \mathbb{R}^d$ are extracted using bilinear interpolation. The local features $\mathbf{f}_i$ from the source images are aggregated into a final feature by computing their weighted mean $\boldsymbol{\mu}$ and variance $\mathbf{v}$ using the barycentric weights.

4

Figure 6: The DTU dataset [1] features various objects captured from front-facing cameras. Our network generalizes well to novel views of novel objects and recovers even highly specular materials like the golden bunny.

The aggregated features are also concatenated with the positional encoding of the ray samples which helps the network to recover high-frequency depth. Hence, the aggregated feature for each ray is defined as:

$$\mathbf{g} = [\boldsymbol{\mu}, \mathbf{v}, \gamma(\mathbf{x})], \tag{2}$$

where $\gamma(.)$ is a positional encoding mapping the position into a higher-dimensional space [7].

Before computing the jump towards the next sample, an important consideration is that a single point sample does not contain sufficient information for an accurate jump prediction. Many real-world scenarios have objects with poor texture or ambiguous depth. If each ray sample independently predicts its own jumps, the final depth map will end up being noisy. We argue here that having knowledge of how the neighbouring rays behave is crucial for resolving ambiguities. Therefore, we add a series of $3 \times 3$ convolutions after the feature aggregation step in order to better constrain the features of each ray.

Finally, the ray features are passed through an LSTM that predicts the displacement $\delta$ along the ray which is used to update our depth $t_{i+1} = t_i + \delta$. This process is iterated a fixed number of times (we use 18 in our experiments) and the final ray sample is considered to be on the surface of the object. This is in stark contrast to NeRF-like models which require samples in the order of hundreds.

Since time consumption increases with the number of ray marching iterations and the number of rays we traverse, we propose to alleviate this problem by employing a coarse-to-fine scheme. Instead of creating rays for each pixel of the target view of size $H \times W$, we first ray march from a down-sampled version at quarter resolution. After several ray marching steps, the computed depth map is upsampled bilinearly to half resolution and the ray marching continues. This process, shown in Fig. 4, iterates until the final full resolution is reached. We observe that this scheme works well since locally-close pixels tend to march together and therefore their depth can be recovered by marching them as a whole. We use three levels of hierarchical depth, each with 10, 5, and 3 ray march steps, respectively.

### 3.3 Color

We obtain per-pixel color by projecting the final ray-marched surface into the three input views and bilinearly sampling both color $\mathbf{c}_i$ and local features $\mathbf{f}_i$ which are concatenated together in $\mathbf{k}_i = [\mathbf{c}_i, \mathbf{f}_i]$. Instead of aggregating $\mathbf{k}_i$ using the barycentric weights, we observed that it is beneficial to allow the network to predict the weights. Therefore, similar to IBRNet [12], we first compute $\boldsymbol{\mu}$ and $\mathbf{v}$ using the barycentric weights in order to capture global information. Afterwards, we concatenate these aggregated features with each per-frame feature vector $\mathbf{k}_i$. Each concatenated feature is fed into a small MLP to integrate both local and global information and predict multi-view aware feature $\mathbf{k}'_i$ and blending weights $w_i \in [0, 1]$. We pool $\mathbf{k}'_i$ into mean and variance by using the weights $w_i$ and map the resulting vector to RGB color using another MLP. We denote the final RGB image with $\tilde{\mathbf{I}}$.

The color loss is computed as the $\ell_1$-loss between the recovered RGB and the ground-truth color. This loss implicitly biases the geometry to lie on the true scene surface since the correct depth produces consistent input view features and the color prediction becomes possible. This allows the network to learn from data sets without ground-truth depth models, but may result in noisy geometry for untextured surfaces due to no direct supervision on the predicted depth.

The reader should further note that we output a full RGB map in one pass of our network. In contrast, NeRF-like methods output a limited number of pixels at a time since their ray-marching step is more expensive and therefore requires to run the network multiple times to complete the full image. This allows our method to use more complex losses like perceptual losses which need to operate on the full image.

### 3.4 Loss with Confidence Estimation

Apart from predicting a correct novel view, it is also valuable to allow the network to predict a confidence for each pixel. This allows to account for possible occlusions. Despite having multiple input views that condition our model, some parts of the target image may not be visible from the working set views and therefore cannot be recreated reasonably. In those cases, the network may hallucinate information based on the context. Yet, it should be aware of this hallucination and therefore assign low confidence to that region.



Figure 5: Results from the Realistic Synthetic dataset [7]. Our network captures both high-frequency detail and view-dependent effects like specular reflections.

In order to predict a confidence map, we draw inspiration from the work of Wagner *et al* [11] which attempt to recover fine-grained explanations from classification networks. The input to their classification network is a pixel-wise blend between the image and a zero image. The loss function attempts to set as many pixels as possible to zero without affecting the classification accuracy. Hence, non-zero pixels are the ones that the network deems important for classification.

In our approach, we choose a similar scheme by defining our loss as a blend between the predicted $\tilde{\mathbf{I}}$ and the ground-truth image $\mathbf{I}$. The blend uses the confidence map $\mathbf{Q}$ which encourages it to be as close as possible to 1 such that most of the pixels are chosen from the predicted image. Then our image loss with confidence estimation is defined as:

$$L = \left\| \mathbf{I} - \left( \tilde{\mathbf{I}} \cdot \mathbf{Q} + \mathbf{I} \cdot (1 - \mathbf{Q}) \right) \right\|_1 + \lambda \left\| 1 - \mathbf{Q} \right\|_2 . \tag{3}$$

## 4 Results

### 4.1 Datasets

We evaluate our method on three datasets. DTU [1] contains real images of various objects and is targeted towards evaluation of MVS methods. We use the train and test splits as defined by pixelNeRF [14]: 88 scenes for training and 15 for testing at a resolution of $400 \times 300$. We use this dataset to test the generalization capabilities of our method. The objects in the test set are different from the ones in the training set, so if the network is able to recover novel views of these novel objects, we can conclude that it learned a general reconstruction method. Results of the generalization to novel objects and novel views can be seen in Fig. 6.

Realistic Synthetic 360° [7] contains synthetic images of objects from the upper hemisphere. The dataset contains eight scenes with images at $800 \times 800$ resolution. The objects exhibit several view-
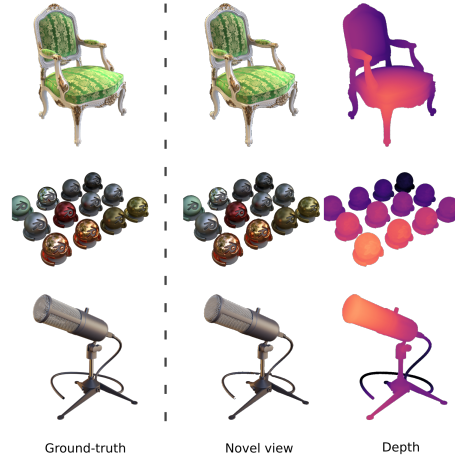
| Method | Setting | DTU [1] | | | Realistic Synthetic 360° [7] | | | Real Forward-Facing [6] | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ |
| pixelNeRF [14] | | 24.14 | 0.887 | 0.224 | 4.36 | 0.46 | 0.44 | 11.266 | 0.388 | 0.757 |
| IBRNet [12] | No per-scene | 25.84 | 0.902 | 0.213 | 19.43 | 0.841 | 0.231 | 16.70 | 0.566 | 0.498 |
| MVSNeRF [2] | optimization | 25.17 | **0.911** | 0.185 | **22.67** | **0.90** | **0.21** | 17.56 | **0.691** | 0.381 |
| Ours | | **26.376** | 0.896 | **0.184** | 20.070 | 0.687 | 0.242 | **18.909** | 0.643 | **0.372** |
| NeRF [7] | | 23.70† | 0.893† | 0.247† | **31.01** | 0.947 | 0.081 | **26.50** | 0.811 | 0.250 |
| MVSNeRF [2] | Per-scene | **29.30** | **0.959** | **0.101** | 27.21 | 0.945 | 0.227 | 26.25 | **0.907** | **0.139** |
| Ours | optimization | 28.093 | 0.913 | 0.165 | 28.425 | **0.952** | **0.070** | 25.206 | 0.803 | 0.218 |

Table 1: Comparison of different methods on multiple real and synthetic datasets. Results with † correspond to NeRF model trained for 9.5h as evaluated by MVSNeRF [2].

| Method | Time | Rays |
|---|---|---|
| NeRF [7] | 6.4 s | 120 k |
| IBRNet [12] | 31 s | 8 k |
| pixelNeRF [14] | 164 s | 300 k |
| Ours | **0.16 s** | **full(640 k)** |

Table 2: Computation time and maximum rays per batch for rendering a $800 \times 800$ image.

dependent effects like specular reflections which must be captured correctly by the network for properly rendering the target view. Results of our network's prediction on this dataset can be seen in Fig. 5.

Real Forward-Facing [6] consists of real images of large scenes scanned with a camera in a forward-facing manner. The dataset contains eight scenes with image size of $1008 \times 756$. We use the train and test split as defined by [12]: every 8th image is selected for testing.

## 4.2 Evaluation

We train our method on the three datasets and distinguish between with and without per-scene optimization. In the case of no scene optimization, we train a generalizable model on the DTU dataset [1] and evaluate on the synthetic [7], a real dataset [6], and the novel scenes from DTU. Tab. **??** shows that our network generalizes to the novel views despite the drastic change in scale and object types. We also train our model with per-scene images similar to NeRF and show that it performs comparable to other generalizable models like MVSNeRF [2] while being significantly faster.

## 4.3 Performance

Previous methods are unable to process the full image at once due to the high computational demand per ray and thus need to run several times with different ray batches to complete the image. We set the ray batches to the maximum size that fits in the memory of an NVIDIA GeForce RTX 3090 and measure the time for rendering a novel view with a resolution of $800 \times 800$ pixels. Tab. 2 shows that our method renders the full image in one forward pass and requires significantly less time than all previous approaches.

## 4.4 Ablation Study

We perform an ablation study of the different components of our network. We train the network on the Lego scene from the synthetic dataset and observe how the network performance is affected.

We first remove the positional encoding from the ray marching.This decreases the depth quality significantly as the network is unable to recover high-frequency details.

Disabling Delaunay for view selection and using a proximity-based method that takes the three closest frames as as the working set shows also a slight decrease in performance.

By default, we use three levels of depth refinement, each with 10, 5, and 3 ray march steps, respectively. We reduce the number of ray march steps to 5, 3, and 1 and observe a slight decrease in

| Setting | PSNR↑ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| No position encoding | 25.692 | 0.899 | 0.082 |
| No Delaunay | 26.764 | 0.919 | 0.075 |
| Fewer ray marches | 27.522 | 0.933 | 0.058 |
| Only $1 \times 1$ convolutions | 26.224 | 0.912 | 0.075 |
| Complete model | **27.918** | **0.937** | **0.056** |

Table 3: Ablation study

performance. The required number of ray marching steps is heavily dependent on the scene complexity. A geometrically simple scene requires few ray marches for the depth to converge to the surface, while more complex ones require more marching steps. The default value we set strikes a good balance between computational cost and accuracy.

Finally, we change the $3 \times 3$ convolutions in the ray marcher to $1 \times 1$ in order to simulate propagating each ray independently with no spatial awareness of the neighbouring rays, similar to other NVS methods. We observe a significant decrease in accuracy, as the rays can no longer leverage spatial information to resolve ambiguities.

## 5   Conclusion

We proposed a network that jointly resolves scene geometry and novel view synthesis from multi-view datasets and is supervised only by image reconstruction loss. We represent the scene geometry as a distance function which we ray march using sphere tracing. Sphere tracing alleviates the memory constraints faced by other methods and allows us to render high resolution images in one forward pass and is thus much faster than previous methods. We further improve the speed by proposing a hierarchical depth refinement which estimates depth in a coarse-to-fine manner.

Finally, we show the generalization capabilities of our network by evaluating on datasets with different scale and object configurations for which we obtain competitive results but with significantly higher frame rates.

## 6   Acknowledgments

## References

[1] Henrik Aanæs, Rasmus Ramsbøl Jensen, George Vogiatzis, Engin Tola, and Anders Bjorholm Dahl. Large-scale data for multiple-view stereopsis. *International Journal of Computer Vision*, 120:153–168, 2016.

[2] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. *arXiv:2103.15595*, 2021.

[3] François Darmon, Bénédicte Bascle, Jean-Clément Devaux, Pascal Monasse, and Mathieu Aubry. Deep multi-view stereo gone wild. *arXiv:2104.15119*, 2021.

[4] Silvano Galliani, Katrin Lasinger, and Konrad Schindler. Massively parallel multiview stereopsis by surface normal diffusion. In *IEEE International Conference on Computer Vision (ICCV)*, 2015.

[5] Michael Goesele, Noah Snavely, Brian Curless, Hugues Hoppe, and Steven M Seitz. Multiview stereo for community photo collections. In *IEEE 11th International Conference on Computer Vision (ICCV)*, 2007.

[6] 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)*, 38(4):29:1–29:14, 2019.

[7] 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 *European Conference on Computer Vision (ECCV)*, 2020.

[8] Andre Rochow, Max Schwarz, Michael Weinmann, and Sven Behnke. FaDIV-Syn: Fast depth-independent view synthesis. *arXiv:2106.13139*, 2021.

[9] Johannes Lutz Schönberger, Enliang Zheng, Marc Pollefeys, and Jan-Michael Frahm. Pixel-wise view selection for unstructured multi-view stereo. In *European Conference on Computer Vision (ECCV)*, 2016.

[10] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3D-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 1119–1130, 2019.

[11] Jörg Wagner, Jan Mathias Köhler, Tobias Gindele, Leon Hetzel, Jakob Thaddäus Wiedemer, and Sven Behnke. Interpretable and fine-grained visual explanations for convolutional neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9097–9107, 2019.

[12] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. IBRNet: Learning multi-view image-based rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.

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

[14] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021.