Addressing the Shape-Radiance Ambiguity in View-Dependent Radiance Fields

Sverker Rasmuson Chalmers University of Technology Gothenburg, Sweden Email: sverker.rasmuson@chalmers.se Erik Sintorn Chalmers University of Technology Gothenburg, Sweden Email: erik.sintorn@chalmers.se Ulf Assarsson Chalmers University of Technology Gothenburg, Sweden Email: uffe@chalmers.se

Abstract—We present a method for handling view-dependent information in radiance fields to help with convergence and quality of 3D reconstruction. Radiance fields with viewdependence suffers from the so called shape-radiance ambiguity, which can lead to incorrect geometry given a high angular resolution of view-dependent colors. We propose the addition of a difference plane in front of each camera, with the purpose of separating view-dependent and Lambertian components during training. We also propose an additional step where we train, but do not store, a low-resolution view-dependent function that helps to isolate the surface if such a separation is proven difficult. These additions have a small impact on performance and memory usage but enables reconstruction of scenes with highly specular components without any other explicit handling of view-dependence such as Spherical Harmonics.

1. Introduction

Recently, NERF [1] and its successors have shown unparalleled results of novel view synthesis of near-photorealistic quality, including view-dependent effects. One pitfall of this approach is the shape-radiance ambiguity, where a case of perfect modeling of view-dependent colors leads to a scenario where every scene can be modelled by, e.g., a sphere, i.e., information about geometry is lost [2]. The model of view-dependence thus has to be good enough to be visually plausible, while not over-fitting the data so that it cannot generalize to new views. While NERF handles most view-dependent effects admirably, there are situations that can be problematic. Storing and computing view-dependent functions can also incur significant overhead to performance and memory.

To alleviate this problem, we propose a method that aims to separate view-dependent and Lambertian colors during reconstruction. In this method, the information in the volume is considered to be purely Lambertian, while the view-dependent components of each camera is captured in a *difference plane* in front of each camera. The colors in the volume cannot be considered Lambertian while semitransparent, but will get closer to this approximation when the volume elements converge to opaque or nearly opaque surfaces. The difference planes store a single value per pixel, which contain information on how much the current pixel deviates from the integrated Lambertian content of each ray sent through the volume (see Figure 1).

The reasoning behind this approach is that if a pixel's color significantly deviates from the Lambertian colors accumulated through the volume, it is by definition viewdependent, and needs to be handled separately.

In some cases it can be hard to separate the colors in purely Lambertian and view-dependent components. We therefore propose the addition of a constraint where a small view-dependent function is computed per voxel. The parameters of this function are not stored, but are only used to increase or decrease the density per voxel to make convergence to the correct surface easier.

A general downside of methods based on neural rendering are the long reconstruction times; usually up to a day or more for a typical scene. Recent papers, such as Plenoxels [3], DirectVoxGo [4] and PERF [5], have shown significant improvements of reconstruction times by solving the optimization problem directly instead of invoking neural networks. To accommodate for view-dependent effects, Plenoxels store Spherical Harmonics for each voxel cell. With 9 extra coefficients per voxel this incur a significant memory cost. DirectVoxGo instead store a grid of features which are evaluated on a shallow MLP. This also requires extra storage and adds a penalty to training and rendering time.

Our method adds low overhead to the PERF-method, and still enables us to handle scenes with highly viewdependent components. This scene can then be used as is or augmented with view-dependent information. The existing geometry should help mitigate the shape-radiance ambiguity in such an endeavor.

2. Related work

Here, we will cover the most relevant previous work, divided into separate categories.

Multi-view Stereo. One classical approach to 3D reconstruction is Multi-View Stereo (MVS), where a set of images are taken of a scene from multiple viewpoints. The images are then paired in sequence to compute depth information about the scene via triangulation. MVS has been well researched over the years, for an overview see Hartley



Figure 1. Results from our algorithm on the drums data set [1]. From left to right: Lambertian color, Lambertian + view-dependent color, view-dependent color only, view-dependent color brightened for easier inspection.

and Zisserman [6] and Seitz et al. [7]. Newer research have to coupled this method with neural networks in Deep Multiview Stereo [8] [9].

Structure from Motion. Another classical approach in Computer Vision is to reconstruct scenes with Structure from Motion (SfM), where image features are used to associate images with another, typically with SIFT features [10] [11] [12]. 3d positions corresponding to image features are jointly optimized with camera parameters to reconstruct the scene. Typically such a reconstruction is sparse, but it can be used as the basis for a dense reconstruction using e.g. PatchMatch Stereo [13] [14].

Novel View Synthesis. Novel View Synthesis is related to 3D reconstruction but focuses on generating novel viewpoints instead of reconstructing the underlying geometry. Well researched areas include image-based rendering and lightfields [15].

Multi-Plane Images (MPIs) is a recent approach where neural networks are used to train semi-transparent images that lie in different planes between the cameras and the scene. These can then be rendered from novel view points by blending the contributions of each such plane intersecting with the view-rays from a virtual camera. [16] [17] [18].

Neural rendering. The interest in neural rendering, where neural networks are coupled with a volume representation of the scene, has exploded in recent years. Neural rendering allows for great reconstruction quality of general scenes, especially when it comes to novel view synthesis, as shown in Neural Volumes [19] and NERF [1]. NERF uses positional encoding for its inputs to the neural network, which give better results for high frequency scene information [20].

Plenty of research has followed this work, for example improving upon the rendering times and view synthesis quality [21] [22] [23] [24]. An overview of this field can be found in the STAR-paper by Tewari et al. [25].

Nerf++ [2] shows that a key component to the ability for NERF to generalize so well to novel views is the structure of the Multi Layer Perceptron (MLP), where the view direction is inserted late in the network. This regularizes the colors to vary smoothly with view direction, and there are less parameters to describe the overall view-dependent color, avoiding over-fitting.

A downside of neural rendering methods are the long times it takes to train the MLP, which can amount to days. Early attempts to speed up training use pretraining or combine it with known methods such as external MVS reconstructions [26] [27] [28].

Recently, Plenoxels [3], DirectVoxGo [4] and PERF [5] have proposed to speed up the reconstruction by direct optimization instead of using neural networks. This improves the reconstruction times by many orders of magnitude compared to the original NERF implementation. View-dependent effects are modeled with Spherical Harmonics in Plenoxels, which incur a significant memory cost. DirectVoxGo use a shallow MLP to account for these effects in a similar way as SNeRG [29], which induces a penalty to training and rendering times.

Our work is built as an addition to the non-linear least squares framework in PERF [5], which has no explicit handling of view-dependent effects. We aim to keep the performance benefits of this framework while still ameliorating the problem of the shape-radiance ambiguity.

3. Shape-Radiance Ambiguity

Granted enough angular resolution of view-dependent colors, almost any shape can satisfy the incoming radiance to a given camera. This is easy to realize by studying Figure 3, where a single (erroneous) point satisfies the incoming radiance for three different views. In reality, they actually correspond to three distinct points on the surface.

In Nerf++ [2], the authors suggest that the structure of the MLP in NERF serves to regularize the view-dependent colors. This in combination with the limited angular resolution available (the size of the view-dependent MLP) works to ameliorate the problem of the shape-radiance ambiguity for novel view synthesis. However, artifacts can still be seen if the goal instead is to recover the reconstructed geometry. In Figure 3, an example of the geometry from a scene reconstructed with NERF is shown, where the radiance field is converted to triangles via marching cubes. The geometry



Figure 2. Illustration of the shape-radiance ambiguity. A single point on a false surface (solid) can satisfy the incoming radiance to the cameras through view-dependent colors, even though these rays correspond to three distinct points on the real surface (dotted).



Figure 3. Geometry from reconstruction with NERF, converted to triangles via marching cubes. Artifacts from the shape-radiance ambiguity can be seen, especially on surfaces where the color is highly view-dependent.

is especially problematic in areas that correspond to highly view-dependent surfaces.

This example shows that additional constraints are needed for NERF and other methods built on viewdependent radiance fields to work as an end-to-end pipeline for 3D reconstruction. In this paper, we propose two lowoverhead mechanisms of separating view-dependent and Lambertian color information for this purpose.

4. Method

As in PERF we represent the volume with a voxel grid, storing color and density per cell. The scene is surrounded with a hierarchical shell of environment maps to help with foreground/background segmentation.

4.1. Volume Rendering

The volume is rendered by, for each pixel in a virtual camera, traversing a ray through the volume and integrating the resulting color according to

$$\mathbf{H} = \sum_{t} T(t)(1 - e^{(-\sigma(t)\delta(t))})\mathbf{c}(t)$$
(1)

for step t, where

$$T(t) = e^{\left(-\sum_{t_n}^t \sigma(t)\delta(t)\right)}$$

is the accumulated throughput, σ is the density, δ is the distance between two adjacent steps, and c is the color.

4.2. Volume Reconstruction

The volume is constructed by invoking Equation 1 for all pixels in the collection of images, and comparing the accumulated color to the reference color at the corresponding pixel. This amounts to minimizing the objective function F with respect to all density and color values $(c_j; \sigma_j)$ such that

$$\min_{c_j;\sigma_j} F,\tag{2}$$

$$F = 0.5 \sum_{i} \sum_{k=1}^{3} (H_{i,k} - r_{i,k})^2$$
(3)

for all pixels i and each color channel k.

4.3. Difference planes

To each camera image we associate a difference plane with identical resolution. For this plane we store one float value per pixel. This value, α_{s_i} , represents how to weight in a view-dependent part of the camera image according to:

$$\mathbf{H}_{i} = (1.0 - e^{(-\alpha_{s_{i}}\sigma_{s})})(\mathbf{r}_{i} - \mathbf{H}_{d_{i}}) + \mathbf{H}_{d_{i}}$$
(4)

where r_i is the reference color in the corresponding pixel i, and

$$\begin{split} \mathbf{H}_{d_i} &= \sum_j e^{-V_j} (1.0 - e^{(-\sigma_j \delta_j)}) \mathbf{c}_j, \\ V_j &= \sum_{k=0}^{j-1} \sigma_k \delta_k \end{split}$$



Figure 4. Illustration of how the difference planes can capture a deviating color at a certain point, in this case induced from a highlight seen only from a specific angle. The deviating camera sees a yellow color (from the light source) in its reference image, compared to most other cameras which see a red diffuse color. The camera can be satisfied by adding green to its corresponding difference plane. The diffuse color in the volume can thus stay red while still satisfying all cameras.

is the Lambertian part of the color accumulated from the ray traversing the volume through voxels j according to Equation 1. The variable α_s is set to zero initially, but can grow given that no acceptable solution can be found using only the Lambertian values in the volume, as part of the minimization of Equation 3. The parameter σ_s can be used as a weight to tune how readily values should be pushed into the specular planes. This value is constant, with a value of $\sigma_s = 0.002$ for all of our experiments.

The value to be minimized, in the same manner as equation 3, is the squared differences of all accumulated values \hat{H}_i and the reference color r_i such that

$$\hat{F} = 0.5 \sum_{i} \sum_{k=1}^{3} (\hat{H}_{i,k} - r_{i,k})^2$$
(5)

for all pixels i and each color channel k.

4.4. Partial derivatives

To allow for this minimization, partial derivatives of equation 5 need to be computed. For all voxel colors c_i we get

$$\frac{\partial F}{\partial \mathbf{c}_j} = \sum_i (\mathbf{H}_i - \mathbf{r}_i) \frac{\partial \mathbf{H}_i}{\partial \mathbf{c}_j}$$

using the chain rule. The last factor can be expanded in the same manner as

$$\frac{\partial \mathbf{H}_{\mathbf{i}}}{\partial \mathbf{c}_{j}} = e^{(-\alpha_{s_{i}}\sigma_{s})} \frac{\partial \mathbf{H}_{d_{i}}}{\partial \mathbf{c}_{j}}$$

This means that the partial derivatives $\frac{\partial \mathbf{H}_d}{\partial \mathbf{c}_j}$ can be computed in the same manner as in PERF. The same derivation also holds for the density values α_j and the corresponding partial derivatives $\frac{\partial \mathbf{H}_d}{\partial \alpha_i}$.

The partial derivative with respect to the new variable α_{s_i} trivially becomes

$$\frac{\partial \mathbf{H}_{\mathbf{i}}}{\partial \alpha_{s_i}} = -\sigma_s e^{(-\alpha_{s_i}\sigma_s)} (\mathbf{r}_i - \mathbf{H}_{d_i}).$$

4.5. Non-highlight colors

To be able to handle view-dependent lighting that does not behave like a highlight, i.e., does not have a clear Lambertian and view-dependent component, we have added an additional step to our method.

NERF and similar methods can capture geometry by utilizing a low-pass filtering function to describe view dependence. If this low-resolution function can find a good approximation of the outgoing radiance for a given point in space, it is more likely that this is a point on a surface than an arbitrary point in the volume.

To emulate this behavior, keeping in mind our goal of performant 3D reconstruction, we define a low-frequency view dependent function for each voxel. This function is implemented as a small environment map of 8x4 pixels. For each voxel, we populate this environment map by blending colors from all cameras, first ray-marching through the volume to account for visibility, see algorithm 1.

This function is then evaluated on how well it succeeded to approximate every given camera color with visibility, by taking the squared error between the reference color in each camera and the color in the environment map for the corresponding direction. For each voxel we get an error Algorithm 1 Algorithm for populating the small viewdependent function for each voxel. Executed once for each voxel.

v = current voxelfor all cameras c do $p = \text{project_voxel}(v, c)$ $T = \text{trace_visibility}(v, c)$ $(\theta, \phi) = \text{map_spherical}(v, c)$ $\text{env_map}(\theta, \phi) += Tp$ $T_{sum} += T$ end for
for $\theta \in [0, \pi], \phi \in [-\pi, \pi]$ do $\text{env_map}(\theta, \phi) /= T_{sum}$ end for

$$E = \frac{1}{\sum_{k} T_{k}} \sum_{k} \sum_{i} T_{k} (c_{i} - p_{i})^{2}$$
(6)

for visibility (to each camera k) T_k , voxel color c_i and predicted color p_i for all cameras k and color channels i.

We have now obtained a cost function for each voxel in the grid, which gives us a measure of how likely it is that a given point in space lies on a surface.

This cost function is added as a weight w to the Cauchy loss following Plenoxels [3] and SnerG [29]:

$$\mathcal{L}_c = \lambda_c \sum_i \log(1 + \lambda_n w \sigma_i^2) \tag{7}$$

for density σ for each voxel *i*, where λ_c and λ_n controls the overall loss and how much the cost function should weigh in respectively. We use the values $\lambda_c = 0.05$ for the artifical scenes and $\lambda_c = 0.01$ for the real scenes, with $\lambda_n = 10$ for both.

This equation is trivially differentiable, and can be implemented in the non-linear least squares framework of PERF by adding a residual per density value σ to be minimized.

With this implementation no expensive view-dependent function needs to be explicitly stored, while the benefits of finding correct geometry are still obtained (see Figure 4.5).



Figure 5. Comparison of a reconstructed detail in the lego scene when using our handling of non-highlight colors (top) compared to a baseline (bottom).

4.6. Auxiliary terms

To ease the computational impact of equation Algorithm 1 and Equation 6, we implement a simple method to skip empty space using the hierarchical formulation in PERF [5]. For each level above the first, we store the voxel grid of the previous level. If the density value for a voxel in that grid is smaller then a threshold $\lambda_v = 1.0$, then that voxel is skipped and treated as if being empty.

For the real scenes from the Tanks and Temples data set [30], we also use a quadratic on the total visibility T for each ray following PERF [5]:

$$\mathcal{L}_s = \lambda_s (-4(T - 0.5)^2 + 1) \tag{8}$$

with $\lambda_s = 0.1$. This term helps to mitigate smoke-like artifacts and keeping the real scenes sparse.

5. Implementation

The difference planes of Equation 4 and Cauchy loss with per-voxel weight of Equation 7 are added to the non-linear least-squares framework used in PERF.

The resolution of the added planes follow the camera resolutions through the hierarchical steps, so that it approximately matches the voxel resolution at any given level of the hierarchy. The values of each new hierarchical level are upsampled using linear interpolation.

The scenes are calibrated using the Structure from Motion software Colmap [14].

Algorithm 1 and Equation 6 are run once in the beginning of each hierarchical level (except the first), and the weights, one float per voxel, are stored in a buffer. The Cauchy loss in Equation 7 is then evaluated in each iteration with these weights.

6. Results

In Figure 6, we test our solution on a number of different 360° scenes from the artificial NERF data set, and from the Tanks and Temples dataset of captured scenes [1] [30]. We use four hierarchical levels with 256^{3} for the hightest level. In most cases, it is possible to separate out view-dependent information, as can be seen in the rendering of the view-dependent parts on the right.

The scenes are reconstructed using a Nvidia GTX 1080 graphics card. In Table 1, we see reconstruction times for when using the additions described in this paper compared to baseline PERF. To make a fair comparison, Equation 7 is added with a constant w = 0.0001 to the baseline case.

For the artificial scenes, we actually have a higher performance due to faster convergence of the reconstruction. For the two real scenes, which have more camera poses than the artificial scenes, the extra overhead is about 20-30%.

7. Discussion

We present a novel way of mitigating the shape-radiance ambiguity in volume reconstruction with view-dependent



Figure 6. Results from our algorithm. From left to right: Lambertian color (from volume), Lambertian + view-dependent color, view-dependent color only, view-dependent color brightened for easier inspection.

colors. Our method requires little extra memory and has a small overall impact on performance, and might even be beneficial for some scenes due to improved convergence. The resulting reconstruction can be used as is with geometry

Timings (seconds)		
	w/ planes	w/o planes
lego	154	194
ficus	107	208
hotdog	267	355
drums	121	273
chair	90	229
material	243	478
ignatius	228	172
family	154	131

TABLE 1. TIMINGS WHEN RECONSTRUCTING WITH OUR TWO ADDITIONS (LEFT COLUMN) AND WITHOUT (RIGHT COLUMN). FOR THE ARTIFICIAL SCENES OUR ADDITION BECOMES A NET SPEED UP, EVEN THOUGH ADDITIONAL COMPUTATION IS PERFORMED. IN THE REAL SCENES (BOTTOM TWO) THE EXTRA COMPUTATION TIME IS MORE NOTICEABLE.

and Lambertian colors, or be further augmented with viewdependent information at a later stage. Given a correct start geometry, a high-resolution function describing viewdependence can be used, without problems with the shaperadiance ambiguity.

7.1. Limitations and Future Work

Our proposed addition is able to successfully capture deviating radiance compared to the Lambertian volume, such as a highlight seen from one or a few cameras. In more complex cases, our small view-dependent function can help with convergence to the correct surface. The next step, which we have not studied in this paper, is to use this separation of Lambertian and view-dependent colors to either explicitly reconstruct the geometry with e.g. marching cubes, and/or to fit a high resolution view-dependent function to describe the incoming radiance with high accuracy.

Acknowledgments

This work was supported by the Swedish Research Council under Grant 2014-4559.

References

- [1] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 12346 LNCS, pp. 405–421, 2020.
- [2] K. Zhang, G. Riegler, N. Snavely, and V. Koltun, "Nerf++: Analyzing and improving neural radiance fields," 2020.
- [3] A. Yu, S. Fridovich-Keil, M. Tancik, Q. Chen, B. Recht, and A. Kanazawa, "Plenoxels: Radiance fields without neural networks," 2021.
- [4] C. Sun, M. Sun, and H.-T. Chen, "Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction," 2021.
- [5] S. Rasmuson, E. Sintorn, and U. Assarsson, "Perf: Performant, explicit radiance fields," 2021.

- [6] R. H. A. Zisserman, "Multiple view geometry in computer vision," 2004.
- [7] S. M. Seitz, B. Curless, J. Diebel, D. Scharstein, and R. Szeliski, "A comparison and evaluation of multi-view stereo reconstruction algorithms," in 2006 IEEE computer society conference on computer vision and pattern recognition (CVPR'06), vol. 1. IEEE, 2006, pp. 519–528.
- [8] Y. Yao, Z. Luo, S. Li, T. Fang, and L. Quan, "Mvsnet: Depth inference for unstructured multi-view stereo," *European Conference* on Computer Vision (ECCV), 2018.
- [9] Y. Yao, Z. Luo, S. Li, T. Shen, T. Fang, and L. Quan, "Recurrent mvsnet for high-resolution multi-view stereo depth inference," *Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [10] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vision*, vol. 60, no. 2, pp. 91– 110, Nov. 2004. [Online]. Available: http://dx.doi.org/10.1023/B: VISI.0000029664.99615.94
- [11] N. Snavely, S. M. Seitz, and R. Szeliski, "Photo tourism: exploring photo collections in 3d," in ACM siggraph 2006 papers, 2006, pp. 835–846.
- [12] O. Özyeşil, V. Voroninski, R. Basri, and A. Singer, "A survey of structure from motion*." *Acta Numerica*, vol. 26, pp. 305–364, 2017.
- [13] M. Bleyer, C. Rhemann, and C. Rother, "Patchmatch stereo-stereo matching with slanted support windows." in *Bmvc*, vol. 11, 2011, pp. 1–11.
- [14] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 4104–4113.
- [15] H. Shum and S. B. Kang, "Review of image-based rendering techniques," in *Visual Communications and Image Processing 2000*, vol. 4067. International Society for Optics and Photonics, 2000, pp. 2–13.
- [16] B. Mildenhall, P. P. Srinivasan, R. Ortiz-Cayon, N. K. Kalantari, R. Ramamoorthi, R. Ng, and A. Kar, "Local light field fusion: Practical view synthesis with prescriptive sampling guidelines," *ACM Trans. Graph.*, vol. 38, no. 4, Jul. 2019. [Online]. Available: https://doi.org/10.1145/3306346.3322980
- [17] T. Zhou, R. Tucker, J. Flynn, G. Fyffe, and N. Snavely, "Stereo magnification: Learning view synthesis using multiplane images," *ACM Trans. Graph.*, vol. 37, no. 4, Jul. 2018. [Online]. Available: https://doi.org/10.1145/3197517.3201323
- [18] M. Broxton, J. Flynn, R. Overbeck, D. Erickson, P. Hedman, M. Duvall, J. Dourgarian, J. Busch, M. Whalen, and P. Debevec, "Immersive light field video with a layered mesh representation," *ACM Trans. Graph.*, vol. 39, no. 4, Jul. 2020. [Online]. Available: https://doi.org/10.1145/3386569.3392485
- [19] S. Lombardi, T. Simon, J. Saragih, G. Schwartz, A. Lehrmann, and Y. Sheikh, "Neural volumes: Learning dynamic renderable volumes from images," *ACM Trans. Graph.*, vol. 38, no. 4, Jul. 2019. [Online]. Available: https://doi.org/10.1145/3306346.3323020
- [20] M. Tancik, P. P. Srinivasan, B. Mildenhall, S. Fridovich-Keil, N. Raghavan, U. Singhal, R. Ramamoorthi, J. T. Barron, and R. Ng, "Fourier features let networks learn high frequency functions in low dimensional domains," *Adv. Neural Inf. Process. Syst.*, vol. 2020-December, pp. 1–24, 2020.
- [21] R. Martin-Brualla, N. Radwan, M. S. M. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections," in *CVPR*, 2021.
- [22] L. Liu, J. Gu, K. Z. Lin, T.-S. Chua, and C. Theobalt, "Neural sparse voxel fields," *NeurIPS*, 2020.
- [23] A. Yu, R. Li, M. Tancik, H. Li, R. Ng, and A. Kanazawa, "PlenOctrees for real-time rendering of neural radiance fields," in *ICCV*, 2021.

- [24] J. T. Barron, B. Mildenhall, M. Tancik, P. Hedman, R. Martin-Brualla, and P. P. Srinivasan, "Mip-nerf: A multiscale representation for antialiasing neural radiance fields," 2021.
- [25] 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.
- [26] M. Tancik, B. Mildenhall, T. Wang, D. Schmidt, P. P. Srinivasan, J. T. Barron, and R. Ng, "Learned Initializations for Optimizing Coordinate-Based Neural Representations," 2020. [Online]. Available: http://arxiv.org/abs/2012.02189
- [27] A. Chen, Z. Xu, F. Zhao, X. Zhang, F. Xiang, J. Yu, and H. Su, "MVSNeRF: Fast Generalizable Radiance Field Reconstruction from Multi-View Stereo," 2021. [Online]. Available: http://arxiv.org/abs/2103.15595
- [28] Y. Liu, S. Peng, L. Liu, Q. Wang, P. Wang, T. Christian, X. Zhou, and W. Wang, "Neural rays for occlusion-aware image-based rendering," arXiv preprint arXiv:2107.13421, 2021.
- [29] P. Hedman, P. P. Srinivasan, B. Mildenhall, J. T. Barron, and P. Debevec, "Baking neural radiance fields for real-time view synthesis," *ICCV*, 2021.
- [30] A. Knapitsch, J. Park, Q.-Y. Zhou, and V. Koltun, "Tanks and temples: Benchmarking large-scale scene reconstruction," ACM Transactions on Graphics, vol. 36, no. 4, 2017.