

# Strata-NeRF : Neural Radiance Fields for Stratified Scenes

Ankit Dhiman<sup>1,2</sup> R Srinath<sup>1</sup> Harsh Rangwani<sup>1</sup> Rishubh Parihar<sup>1</sup>  
 Lokesh R Boregowda<sup>2</sup> Srinath Sridhar<sup>3</sup> R Venkatesh Babu<sup>1</sup>

<sup>1</sup>Vision and AI Lab, IISc Bangalore <sup>2</sup>Samsung R & D Institute India - Bangalore <sup>3</sup>Brown University

## Abstract

Neural Radiance Field (NeRF) approaches learn the underlying 3D representation of a scene and generate photo-realistic novel views with high fidelity. However, most proposed settings concentrate on modelling a single object or a single level of a scene. However, in the real world, we may capture a scene at multiple levels, resulting in a layered capture. For example, tourists usually capture a monument’s exterior structure before capturing the inner structure. Modelling such scenes in 3D with seamless switching between levels can drastically improve immersive experiences. However, most existing techniques struggle in modelling such scenes. We propose Strata-NeRF, a single neural radiance field that implicitly captures a scene with multiple levels. Strata-NeRF achieves this by conditioning the NeRFs on Vector Quantized (VQ) latent representations which allow sudden changes in scene structure. We evaluate the effectiveness of our approach in multi-layered synthetic dataset comprising diverse scenes and then further validate its generalization on the real-world RealEstate10K dataset. We find that Strata-NeRF effectively captures stratified scenes, minimizes artifacts, and synthesizes high-fidelity views compared to existing approaches. <https://ankitatiisc.github.io/Strata-NeRF/>

## 1. Introduction

Novel view synthesis is an ill-posed problem widely encountered in various areas such as augmented reality [28, 32], virtual reality [13], etc. A paradigm change for solving these kinds of problems was brought by the introduction of Neural Radiance Fields (NeRF) [38]. NeRFs are neural networks that take in the spatial coordinates and camera parameters as input and output the corresponding radiance field. Earlier version of NeRFs enable the generation of high-fidelity novel views for bounded scenes, significantly improving over existing techniques like Structure From Motion [52]. Further, the capability of NeRFs have been recently extended to model unbounded scenes by Mip-NeRF 360 [2]. This enabled NeRFs to model complex real-world

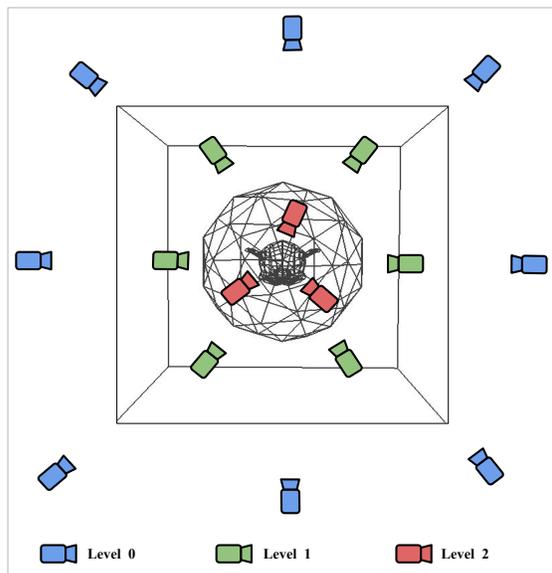


Figure 1. Top, wireframe view of a multi-layered stratified scene with three levels (monkey head inside sphere inside a cube). The camera colors indicate views of a specific level. Strata-NeRF enables high-quality reconstruction of such stratified scenes using a single neural network.

scenes, where the scene content can exist at any distance from the camera.

However, similar to unboundedness in scenes, hierarchies in scenes are also natural. For example, images captured in a house can be categorized into images captured outside and inside across various rooms. Modelling such hierarchical scenes jointly for all levels through a NeRF could be particularly useful in cases of Virtual Reality applications. As it would not require switching to a different NeRF for each level, reducing memory requirement and latency in switching. Further, as the different hierarchies of a scene usually share texture and architectural commonalities, it could lead to effective knowledge sharing and reduce the requirement of training independent models. For tackling the above novel objective, we introduce a paradigm of scenes that can be deconstructed into several tiers, termed “Stratified Scenes”. A “stratified” scene has several levels or

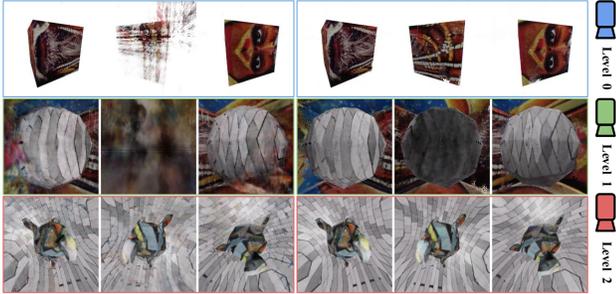


Figure 2. Novel views for stratified scene in Figure 1, from Mip-NeRF 360 [2] (left) and our method “*Strata-NeRF*” (right). Existing methods struggle to capture stratified scenes with a single network while ours produces sharp results.

groupings of structure (Figure 1). In our work, we first propose a synthetic dataset of stratified scenes, i.e. scenes having multiple levels. This dataset comprises scenes from two categories: (i) Simpler geometry, such as spheres, cubes, or tetrahedron meshes, and (ii) Complex geometry, which closely emulates a real-world setup.

On such datasets, we find methods such as Mip-NeRF 360 perform well for each level of the hierarchy independently, but produce unsatisfactory results when images from all hierarchical levels are used together for training (Figure 3). This can be attributed to the continuous nature of NeRFs, which is unsuitable for modelling the sudden changes in scenes with shifts in hierarchical levels. Hence, in this work, we introduce *Strata-NeRF* that explicitly aims to model the hierarchies by conditioning [26, 42, 43, 72, 48] the NeRF on Vector Quantized (VQ) latents. The VQ latents enable the modelling of discontinuities and sudden changes in the scene, as they are discrete and less correlated with others [62]. In practice, the VQ conditioning is achieved by introducing two lightweight modules: the “Latent Generator” module that compresses the implicit information in encoded 3D positions to generate VQ latent code, which is directed through the “Latent Routing” module to condition various layers of radiance field. The additional parameters introduced through these modules are significantly less than training an independent NeRF model for each level, leading to a significant reduction in memory.

For evaluating the proposed *Strata-NeRF* we first test on the proposed synthetic *Stratified Scenes* dataset, where we find that *Strata-NeRF* learns the structure in scenes across all levels. In contrast, other baselines produce cloudy and sub-optimal novel views (Figure 2). Further, to test the generalizability of the proposed method on real-world scenes, we utilize the high-resolution RealEstate10K dataset. We find that *Strata-NeRF* significantly outperforms other baselines and produces high-fidelity novel views without artifacts compared to baselines. This is also observed quantitatively through improvement in metrics, where it establishes

a new state-of-the-art. In summary,

- We first introduce the task of implicit representation for 3D stratified (hierarchical) scenes using a single radiance field network. For this, we introduce a novel synthetic dataset comprising of scenes ranging from simple to complex geometries.
- For implicit modelling of the stratified scenes, we propose *Strata-NeRF*, which conditions the radiance field based on discrete Vector-Quantized (VQ) latents to model the sudden changes in scenes due to change in hierarchical level (i.e. strata).
- *Strata-NeRF* significantly outperforms the baselines across the synthetic dataset and generalizes well on the real-world scene dataset of RealState10k.

## 2. Related Work

Generating photo-realistic novel views from densely sampled images is a classical problem. Earlier methods solved this issue using light-field-based interpolation techniques [12, 21, 31]. These techniques interpreted the input images as 2D slices of a 4D function - the light field. The only caveat in these methods is their overreliance on dense views. Another popular technique is Structure From Motion (SFM) which reconstructs 3D structure of a scene or an object by using a sequence of 2D images. We suggest readers to read survey papers [52, 41] to understand SFM methods in detail. Shum *et al.* [54] also provides an excellent review on traditional image based rendering techniques.

**Neural Volume Reconstruction.** NeRF [38] has shown remarkable results in encoding the 3D geometry of a scene implicitly using the multi-layer perceptron (MLP). Specifically, it trains an MLP, which takes 3D position and a viewing direction to predict colour and occupancy. Many papers have extended this idea to solve different scenarios such as dynamic scenes, low-light scenes, synthesis from fewer views, accelerating the performance etc. Mip-NeRF [1] mitigates the problem of aliasing when a novel view is generated at a different resolution. MVNeRF [9] generalizes across all the scenes and optimizes the geometry and radiance field using only a few views. NerfingMVS [67] utilizes conventional SFM reconstruction and learning-based priors to predict the radiance field. UNISURF [40] combines implicit surface models and radiance fields to render both surface and volume rendering.

AR-NeRF [28] replaced pin-hole based camera ray-tracing with aperture camera based ray-tracing. DiVeR [68] uses a voxel based representation to learn the radiance field, Mip-NeRF 360 [2] improves view synthesis on the unbounded scenes and also proposed an online distillation scheme which significantly reduced the training and inference time. Neural Rays [35] solves the occlusion problem by predicting the visibility of the 3D points in their

representation. Scene Representation Transformers [51] uses Vision Transformers [15] to infer latent representations to render the novel views. Further, many methods [34, 20, 49, 71, 56, 25, 64] have been proposed to improve the slow training and inference time for neural radiance field based methods. Despite many works, no work has focused on modelling the *stratified* scenes.

**NeRF Extensions.** Relighting discusses how to model different types of light and then using this model to re-light a scene [36, 3, 55, 63, 23]. Breaking the myth that radiance field can only be used in small and bounded scenes, recent methods [57, 61, 50] have scaled it to large-scale city scenes. Another line of work focuses on modelling the dynamic scenes with presence of moving objects [42, 69, 33, 45, 16, 60, 19, 43] through NeRFs.

**Neural Radiance Fields and Latents.** Recently, a lot of methods have made use of the latents to bring generative capabilities to neural radiance fields. GRAF [53] uses disentangled shape and appearance latent codes to generalize on an object category. For viewpoint invariance, they used typical GAN based training. Pi-GAN [7] uses volumetric rendering equations for consistent 3D views in a generative framework. Pixel-NeRF [72] learns a scene prior to generalize across different scenes. GSN [14] decomposes the radiance field of a scene into local radiance fields by conditioning on a 2D grid of latent codes. Code-NeRF [26] learns the variation of object shapes and textures across by learning separate latent embeddings. LOLNeRF [48] uses a shared latent space which conditions a neural radiance field to model shape and appearance of a single class. PixNeRF [6] extends Pi-GAN [7] and maps images to a latent manifold allowing object-centric novel views given a single image of an object. NeRF-W [37] optimizes latent codes to model the scene variations to produce temporally consistent novel view renderings. In contrast to these methods, we propose conditioning NeRF on learnable Vector Quantized latents.

**Vector Quantized Variational Autoencoders (VQ-VAE) [62]:** VQ-VAE uses vector quantization to represent a discrete latent distribution. VQ-VAE has shown applications in Image Generation [46, 44], speech and audio processing [22, 65]. Further, its extension like VQ-VAE2 [46] uses hierarchical latent space for high-quality generation.

### 3. Preliminaries

NeRF represents a scene as an implicit function  $f : (X, d) \rightarrow (c, \sigma)$  which maps a 3D position  $X = (x, y, z)$  and  $d = (\theta, \phi)$  to a color  $c = (r, g, b)$  and occupancy density  $\sigma$ . An MLP parametrizes this implicit function  $f$ . Before sending the inputs  $X$  and  $d$  through the network, a positional encoding is used to project them in a high dimensional space [58]. Finally, the volume rendering [27] procedure enables NeRF to represent scenes with photo-realistic

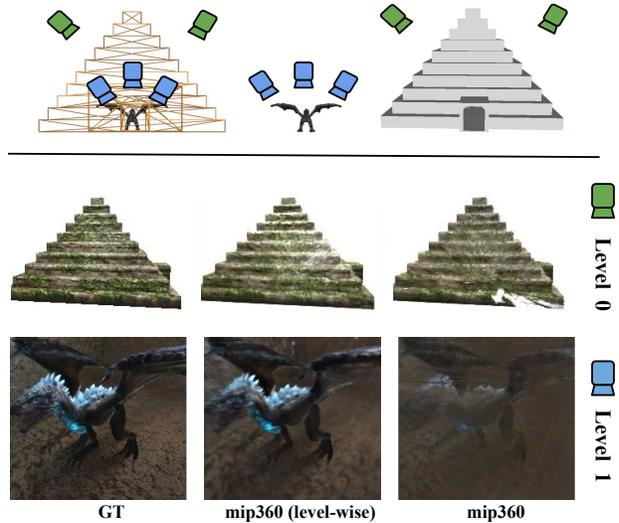


Figure 3. Analysis on “Dragon in pyramid” scene. The top row shows the layout of the levels in 3D scene. Observe that baseline works fine on the scenes when trained individually. Artefacts occur when the baseline is trained on views from the entire scene.

rendering from novel camera viewpoints.

**Volume Rendering.** At the crux of NeRF lies the volume rendering equation. A ray  $r(t) = o + td$  is cast from the camera center  $o$  through the pixel along direction  $d$ . The pixel’s color value is estimated by integrating along the ray  $r(t)$  as described in Eq. 1

$$c(r) = \int_{t_n}^{t_f} T(t)\sigma(r(t))c(r(t), d) dt \quad (1)$$

where transmittance  $T(t) = \exp(-\int_{t_n}^t \sigma(r(s)) ds)$  is the probability that a ray passes unhindered from the near plane ( $t_n$ ) to plane ( $t$ ) and use this probability to integrate till far plane ( $t_f$ ). In Mip-NeRF [1], a ray  $r(t)$  is divided into intervals  $T_i = [t_i, t_{i+1})$  which corresponds to a conical frustum. For each interval  $T_i$ , it computes the mean and variance  $(\mu, \Sigma)$  and uses it for integrated position encoding as illustrated in Eq. 2.

$$\gamma(\mu, \Sigma) = \left\{ \begin{array}{l} \left[ \begin{array}{l} \sin(2^l \mu) \exp(-2^{2l-1} \text{diag}(\Sigma)) \\ \cos(2^l \mu) \exp(-2^{2l-1} \text{diag}(\Sigma)) \end{array} \right]_0^{L-1} \end{array} \right. \quad (2)$$

This solves the aliasing issue in the original NeRF. Mip-NeRF 360 [2] proposed coarse-to-fine online distillation for proposal sampling, which efficiently reduces the training time as the proposed MLP only predicts density. They also proposed ray parametrization and regularisation techniques to alleviate hanging artifacts in unbounded scenes. We’ll refer Mip-NeRF 360 [2] as *mip360* in all our discussions. We choose *mip360* [2] as the baseline for all our experiments.

Table 1. A quantitative comparison of mip360 (level-wise) and mip360 (all views) on “Dragon in pyramid” scene.

	Level 0			Level 1		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
mip360 (level-wise)	31.5390	0.9181	0.1304	29.8560	0.8133	0.3484
mip360	30.8847	0.9006	0.1367	24.3876	0.7055	0.5163

## 4. Motivation

The majority of real-world scenarios are stratified with multiple levels. For example, a commodity store has exterior and interior structures. This work addresses an essential question for such stratified scenes: *Can a single radiance field learn such hierarchical scenes?* This section introduces and discusses our observations on one such stratified scene: “Dragon in Pyramid”, as illustrated in Figure 3. The outer structure of “Dragon in Pyramid” is a Mayan pyramid that has a dragon inside it. To validate our claim, we first train the baseline model on each level, i.e., on outer pyramid views and inner views (focusing dragon) independently. We refer to these separately trained models as *mip360 (level-wise)*. Then, we train a single *mip360* model using the outer and inner views for the scene. The term “level” in our work refers to each level in a stratified scene. In the scene depicted in Figure 3, level 0 denotes the pyramid’s outer construction, while level 1 denotes the pyramid’s interior structure, which contains a dragon.

Table 1 shows that the baseline model performs remarkably well when trained separately on each level. In comparison, the metric values for the baseline model trained jointly on both levels of stratified scene declines. PSNR at level 1 is 24.39 dB, a 5.47 dB reduction compared to mip360 (level-wise). Similarly, performance in level 0 has declined, but less dramatically than in the inner level. This pattern is observed across all metrics. Furthermore, the qualitative results illustrated in Figure 3 backs up the quantitative study’s findings. Figure 3 indicates that mip360 (level-wise) generates novel views on par with the ground truth. However, shown in Figure 3, the jointly trained model has white artifacts on the pyramid’s outer structure and haziness in front of the dragon inside the pyramid. This demonstrates that current radiance field networks have issues while learning a 3D representation of a stratified scene. We perform a similar experiment for a RealEstate10K scene in Appendix E.1 in the supplementary material.

## 5. Method

This section describes our method : *Strata-NeRF* for stratified scenes. We generate latent codes with the latent generator described in Section 5.1. This latent code is fed into the radiance field architecture through the latent router, described in Section 5.2. Figure 4 depicts the overall architecture of Strata-NeRF. We adopt the base neural radiance

field architecture proposed in mip360 [2].

### 5.1. Latent Generator

A latent space reflects the scene’s “compressed” representation. It has been shown in various works that this space has rich properties. VQ-VAE [62] learns a codebook to model the discrete distribution of the latent space of a variational-autoencoder. The encoder’s output is compared to all of the vectors in the codebook. The nearest vector is fed into the decoder as input. Since most data in the world is discrete, VQ based models have been highly successful in image generation [17], speech encoding [62], and other applications. In a stratified scene, the definition of level is also discrete. Hence, our method employs VQ-VAE as a latent generator because of their proven success in representing discrete distributions.

We use Integrated Positional Encoded (IPE) [2]  $\gamma(\mathbf{x})$  as input to our latent generator. We encode  $\gamma(\mathbf{x})$  and then search the codebook for the closest vector. After that, the closest vector from the codebook is used to condition the radiance field network. Specifically,  $\gamma(\mathbf{x})$  is passed through a set of two hidden layers to generate an encoded input  $\mathbf{z}$ . The encoded latent code  $\mathbf{z}$  is then passed through the quantizer bottleneck to determine the quantized latent code  $\mathbf{z}_e$ , where  $\mathbf{z}_e \in E$ ; where  $E \in R^{N \times D}$  is the codebook;  $N$  is the number of vectors in the codebook, and  $D$  is the dimension of the latent space.  $\mathbf{z}_e$  is then supplied into the decoder network, which consists of two hidden layers, to yield  $\mathbf{y}$  as the reconstructed output of  $\gamma(\mathbf{x})$ . The quantized latent  $\mathbf{z}_e$  is also sent into the radiance field network through the “Latent Router” block. Loss for this variational autoencoder (VAE) block is defined as follows:

$$\mathcal{L}_{vq} = \|\gamma(\mathbf{x}) - \mathbf{y}\|_2^2 + \|\text{sg}(\mathbf{z}_e) - \mathbf{z}\|_2^2 + \beta \|\mathbf{z}_e - \text{sg}(\mathbf{z})\|_2^2 \quad (3)$$

The “Latent Generator” module based on VAE is jointly trained with the NeRF through backpropagation.

### 5.2. Latent Router

The Latent Router block is inspired by the CodeNeRF architecture [26], in which shape and texture latent codes are sent to the NeRF MLP through a residual connection. In our architecture, the quantized latent codes  $z_e$  that are generated in the “Latent Generator” block are input to the Radiance field after passing through an MLP layer in the Latent Router as shown in Figure 4.

### 5.3. Training Strata-NeRF

For training Strata-NeRF, we utilize the losses suggested by mip360 [2] as we use a similar radiance field design.  $\mathcal{L}_{recon}(c(r, t), c^*(r))$  denotes the reconstruction loss between the estimated colour along a ray and the actual colour value.  $\mathcal{L}_{dist}(s, w)$  is the distortion loss where  $s$  is the normalized ray distances and  $w$  is the weight vector. Note that

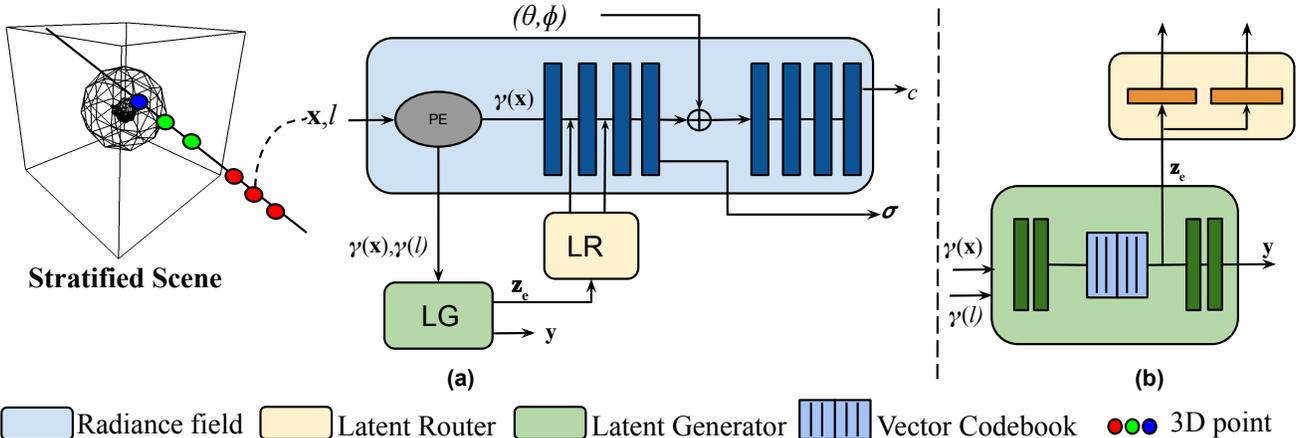


Figure 4. For each 3D point along the projected ray, we generate a latent code using our “Latent generator” module. The generated latent code is routed to the MLP using “Latent Router”. Vector Codebooks learn the discrete distribution of positionally encoded 3D points. (a) Our model’s end-to-end architecture; (b) components of the “Latent Generator” and “Latent Router” blocks.

Table 2. Characteristic Comparison of the proposed methods

Method	Discrete Representation	Photometric Losses	VAE loss
NeRF [38]	✗	✓	✗
mip360 [2]	✗	✓	✗
Plenoxel[70]	✓	✓	✗
Instant-NGP[39]	✓	✓	✗
TensoRF[8]	✓	✓	✗
Ours	✓	✓	✓

we don’t alter anything in the proposal MLP. More details are provided in mip360 [2]. The total loss for Strata-NeRF is given as:

$$\mathcal{L}_{total} = \mathcal{L}_{recon}(c(r, t), c^*(r)) + \lambda_1 \mathcal{L}_{dist}(s, w) + \lambda_2 \mathcal{L}_{vq} \quad (4)$$

We use  $\lambda_1 = 0.01$ ,  $\lambda_2 = 0.1$  and  $\beta = 1.0$  across all our experiments, as they work robustly [2] for *Strata-NeRF*.

## 6. Experiments

We discuss implementation details in Section 6.1. Section 6.2 discusses the dataset used for evaluating our method with other baselines. In Section 6.3, we present quantitative and qualitative comparison with the baseline methods. Additionally, we discuss the ablations for the proposed method.

### 6.1. Implementation Details

Our method builds on mip360 [2] as the base radiance field. We use a latent generator network which consists of an encoder-decoder architecture and a vector-codebook. The encoder has two linear layers of hidden size 48, and the decoder has one linear layer of hidden size 96. The output dimension of our decoder matches the output from Integrated Positional Encoding (IPE) block. The size of

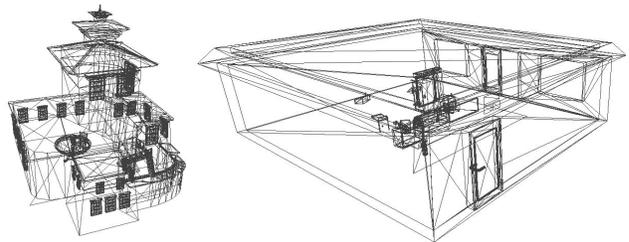


Figure 5. Skeleton mesh of the stratified scenes : Bhutanese House and Coffee Shop. More details are in the supplementary material.

our codebook is 1024, and the dimension of each vector in the codebook is 48. We condition the neural radiance field through the latent generated after the quantization step in the latent generator. We use a Latent routing module consisting of two linear layers of hidden-size 256. As illustrated in Figure 4, the output of the linear layer in the routing module conditions the first two layers of the radiance field network. We employ the losses outlined in Section 5. On each scene, we train our approach for 150k iterations. We use Adam [29] optimizer with a learning rate of  $1e^{-6}$ . Further details are provided in supplementary material.

### 6.2. Evaluation Dataset

Most of the radiance field methods evaluate their results on the synthetic (Blender) and real-world (LLFF) datasets proposed in NeRF [38]. These scenes either include a solitary object on a white background or a frontal view of a natural scene. According to our description of stratified scenes, these datasets has only one level. Even large-scale reconstruction datasets like TanksandTemples [30] are not representative of our setting as they only have views either inside or outside of the structure. Similarly, Scannet [11] a dataset for real-world interior scenes, lacks the characteristics of

Table 3. Quantitative evaluation on test-set against baselines discussed in Section 6.1. Each column is depicts the **best** and **second best**.

	Cube-Sphere-Monkey											
	Level 0			Level 1			Level2			Total		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Nerf [38]	28.3314	0.9383	0.1034	18.1806	0.4976	0.4981	22.1178	0.5995	0.3825	22.8766	0.6784	0.3280
mip360 [2]	28.3149	0.9298	0.1156	19.0443	0.5343	0.4930	24.9136	0.7326	0.3245	24.0909	0.7322	0.3110
Plenoxels [70]	25.3547	0.9169	0.1238	13.1148	0.3320	0.6895	21.5568	0.6523	0.3803	20.0087	0.6337	0.3979
Instant-NGP [39]	28.2104	0.9168	0.1123	14.3648	0.1830	0.7216	17.6914	0.2744	0.5997	20.0889	0.4581	0.4779
TensorRF [8]	<b>32.0077</b>	<b>0.9532</b>	<b>0.0692</b>	13.7487	0.1537	0.7106	13.0075	0.2496	0.6886	19.5880	0.4521	0.4894
Ours	26.9335	0.9298	0.1255	<b>25.7088</b>	<b>0.7738</b>	<b>0.2959</b>	<b>26.1912</b>	<b>0.8172</b>	<b>0.2549</b>	<b>26.2778</b>	<b>0.8403</b>	<b>0.2254</b>

	Bhutanese House											
	Level 0			Level 1			Level 2			Total		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Nerf [38]	11.4478	0.6917	0.3711	17.1209	0.5886	0.7078	18.3918	0.6952	0.6591	15.6535	0.6585	0.5793
mip360 [2]	26.6240	0.9002	0.2062	24.5946	0.7296	0.4739	29.4225	<b>0.8577</b>	0.4156	26.8804	0.8291	0.3652
Plenoxels [70]	15.2205	0.7752	0.3052	13.0386	0.4670	0.6703	19.3050	0.5819	0.5886	15.8547	0.6080	0.5214
Instant-NGP [39]	23.9791	<b>0.9217</b>	<b>0.1500</b>	24.7316	0.7009	0.4237	27.6617	0.8136	<b>0.3786</b>	25.4575	0.8121	<b>0.3174</b>
TensorRF [8]	13.8880	0.7607	0.3142	17.0244	0.4856	0.6421	16.8170	0.6306	0.6332	15.9098	0.6256	0.5298
Ours	<b>27.6842</b>	0.9046	0.2045	<b>24.9180</b>	<b>0.7371</b>	<b>0.4616</b>	<b>29.4646</b>	0.8575	0.4172	<b>27.3556</b>	<b>0.8331</b>	0.3611

Table 4. Quantitative evaluation on test-set against baselines discussed in Section 6.1. Each column is depicts the **best** and **second best**.

	Coffee Shop											
	Level 0			Level 1			Level2			Total		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Nerf [38]	06.7446	0.6197	0.4698	16.1398	0.4915	0.7982	12.8889	0.4213	0.8158	11.9244	0.5108	0.6946
mip360 [2]	26.2073	0.8825	<b>0.1867</b>	27.0500	0.8086	0.3785	<b>34.2023</b>	<b>0.9362</b>	<b>0.1950</b>	29.1532	0.8757	0.2534
Plenoxels [70]	19.3204	0.7968	0.2579	12.3871	0.4044	0.6904	22.4325	0.6856	0.4585	18.0467	0.6289	0.4689
Instant-NGP [39]	29.9425	0.9324	0.0992	28.1040	0.8193	0.3452	29.6574	0.8680	0.2621	29.2347	0.8732	<b>0.2355</b>
TensorRF [8]	<b>33.0337</b>	<b>0.9435</b>	<b>0.0692</b>	19.3115	0.5331	0.6580	21.1852	0.7169	0.4594	24.5102	0.7312	0.3955
Ours	26.4499	0.8802	0.1939	<b>28.6392</b>	<b>0.8403</b>	<b>0.3450</b>	33.2692	0.9254	0.2243	<b>29.4528</b>	<b>0.8819</b>	0.2544

	Dragon In Pyramid											
	Level 0			Level 1			Level 2			Total		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
Nerf [38]	14.6405	0.6595	0.3800	20.8368	0.6052	0.6856	-	-	-	17.7386	0.6323	0.5328
mip360 [2]	30.8758	0.9006	0.1367	24.3890	0.7054	0.5163	-	-	-	27.6324	0.8030	0.3265
Plenoxels [70]	13.0667	0.6247	0.4217	14.5126	0.3572	0.6498	-	-	-	13.7896	0.4910	0.5358
Instant-NGP [39]	23.9054	0.9010	0.0949	24.7389	0.6594	0.4664	-	-	-	24.3222	0.7802	0.2807
TensorRF [8]	<b>35.3015</b>	<b>0.9632</b>	<b>0.0414</b>	19.5573	0.5221	0.6809	-	-	-	27.4294	0.7427	0.3611
Ours	29.4773	0.8700	0.1699	<b>26.1722</b>	<b>0.7489</b>	<b>0.4573</b>	-	-	-	<b>27.8248</b>	<b>0.8095</b>	<b>0.3136</b>

a stratified dataset. Because of the direct unavailability of stratified scenes, we built our own dataset that replicates the intended “stratified” scenario. We create a synthetic scene dataset using a mesh-editing software Blender [10] and real scene dataset by altering RealEstate10K dataset which was proposed for the camera localization task.

The proposed synthetic dataset has two important variations based on: (a) the number of stratified levels and (b) the geometric complexity. We classify based on the geometry’s complexity as follows: (a) *Simple Scenes*: Stratified scenes using geometric components such as the sphere, cube, and so on; and (b) *Complex Scenes*: Stratified scenes that mimic real-world scenes. For Simple Scenes, we leverage models and textures provided by Blender [10]. We utilized publicly available graphical models and composited them to create a real-world configuration for Complex scenes. For example, to design the “Coffee shop” scene, we selected a building structure for the outer level and walls and glasses for the intermediate level structure. For the core level, we com-

posited elements such as a cash register, coffee cups, and so on to simulate a real-world coffee-shop scene. To avoid photo-metric changes, we use fixed illumination. For each stratified level, the camera settings : field of vision and focal length are fixed. Each scene is rendered at  $200 \times 200$  resolution. The camera viewpoint are sampled evenly from the curved surface of a hemisphere and then randomly divided into train, validation, and test sets. Inner objects in *Simple Scenes* are rendered from the surface of a sphere. Figure 5 depicts the proposed dataset’s skeletal meshes. Further information on dataset is present in Appendix B in the supplementary material.

**RealEstate10K dataset.** We extracted four scenes “Spanish Colonial Retreat in Scottsdale Arizona”, “139 Barton Avenue Toronto Ontario”, “31 Brian Dr Rochester NY” and “7 Rutledge Ave Highland Mills” from RealEstate10K dataset. We manually inspected and removed regions which had dynamic components in them. More details about converting RealEstate10k dataset for our stratified setting is

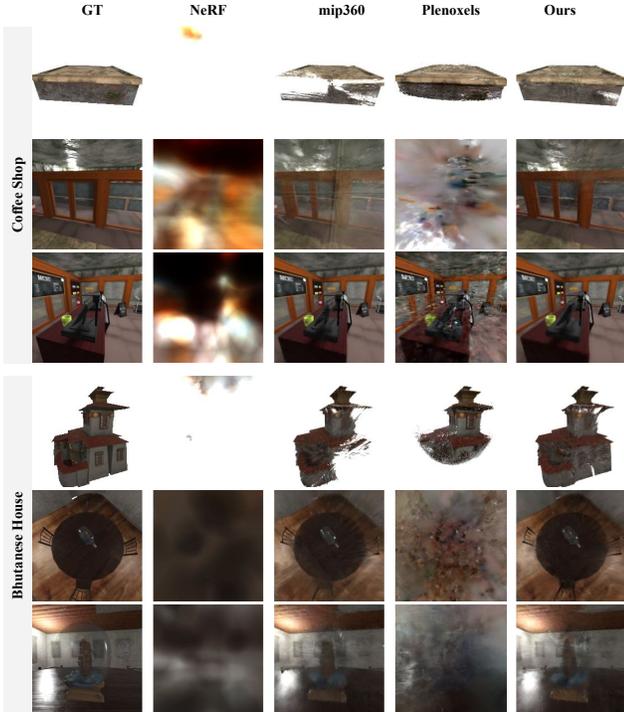


Figure 6. (From top to bottom) Qualitative results on the proposed synthetic datasets (Figure 5). Each row represents a novel view from a level of the stratified scene. The ground-truth (GT) is shown in Column 1. Compared to baselines (Column 2-4), our method’s (Column 5) renderings are more consistent to GT.

provided in Appendix C the supplementary material.

### 6.3. Evaluation

We present quantitative and qualitative analysis of *Strata-NeRF* on the datasets described in Section 6.2.

**Baselines.** We compare our model with NeRF [38], mip360 [2], Instant-NGP [39], TensorRF [8] and Plenoxels [70]. We chose Plenoxels [70] for comparison because it uses sparse-voxel representation which already discretizes the continuous 3D space, which can be useful in stratified scenes. It is worth noting that the sizes of the synthetic scenes in our dataset differ. As a consequence, the authors’ recommended configuration file did not produce the optimal results. As a result, we modified the configuration files for unbounded scenes released by the creators of mip360 [2] to improve performance. For Instant-NGP [39], TensorRF [8] and Plenoxels [70], we change the hyperparameters like bound and scale as suggested in the official implementations. More information is in Appendix D in the supplementary material. Table 2 provides an overview of baselines.

**Quantitative Results.** Table 3 & 4 shows the average PSNR, SSIM [66] and LPIPS [73] for each stratified level in unseen test views. We find that our method surpasses

Table 5. Quantitative comparison of our model and baseline on “139 Barton Avenue” scene of RealEstate10K dataset.

	Metrics	Level 0	Level 1	Level 2	Level 3	Level 4	Level 5
mip360 [2]	PSNR $\uparrow$	18.086	16.496	24.459	20.862	17.479	10.999
	SSIM $\uparrow$	0.618	0.595	0.771	0.702	0.584	0.409
Ours	PSNR $\uparrow$	<b>23.164</b>	<b>21.665</b>	<b>25.236</b>	<b>24.156</b>	<b>22.879</b>	<b>25.409</b>
	SSIM $\uparrow$	<b>0.826</b>	<b>0.757</b>	<b>0.789</b>	<b>0.791</b>	<b>0.753</b>	<b>0.782</b>

Table 6. Quantitative comparison of our model and mip360 baseline on Six Layer Scene.

Dataset	Levels	mip360 [2]	Ours	mip360 [2]	Ours
<i>Spanish Colonial Retreat</i>	5	20.106	<b>22.514</b>	0.622	<b>0.685</b>
<i>31 Brian Dr Rochester</i>	4	23.273	<b>28.026</b>	0.715	<b>0.835</b>
<i>139 Barton Avenue</i>	6	18.991	<b>23.433</b>	0.642	<b>0.780</b>
<i>7 Rutledge Ave</i>	7	19.621	<b>25.040</b>	0.566	<b>0.791</b>

other methods across all metrics most of the time. The baseline mip360 [2] works fine for the exterior structure but fails for the inner layers in the “Cube-Sphere-Monkey” scene. *Strata-NeRF*, on the other hand, offers superior metrics at all stratified levels. The baseline models do well in the outer scene but perform sub-optimally in the inner levels, especially in level 1. These outcomes demonstrate that our method outperforms the baseline models significantly.

Table 6 shows the summary of average PSNR and SSIM for all the levels in a scene for RealEstate10K dataset. In this case, we only compare our method with mip360 as it is the best performing one among others on the synthetic dataset. We observe that our method outperforms the baseline method in all scenarios. Further, we present level-wise result for a specific scene in Table 5. We observe that for real datasets *with increase in number of levels, the magnitude of performance improvement increases*, which demonstrates the effectiveness of the proposed approach. Further, we also compare Instant-NGP [39] and TensorRF [8] on a RealEstate10K scene in Appendix E.2 in the supplementary material.

**Qualitative Results.** Figure 6 & 8 depicts the qualitative results for the synthetic dataset scenes described in Section 6.2. We observe that NeRF [38] performs poorly regardless in majority of scenarios. The generated novel views for “Coffee Shop” are poor. It only works well in level 0 of “Cube-Sphere-Monkey” dataset. mip360 [2] outperforms NeRF but falters in level 1. Furthermore, in level 0 of the “Cube-Sphere-Monkey” dataset, mip360 only generates a white patch with no visible structure. For RealEstate10K dataset, it can be observed in Figure 7 that mip360 generates blurry results compared to our approach. Further, we find that our approach generates consistent and structurally salient novel views throughout all levels and scenes. We show qualitative results for Instant-NGP [39] and TensorRF [8] in Appendix E.2 in the supplementary material.

**Worst Case Analysis.** When comparing different methods, average metrics are often insufficient to determine which method is superior to the others. As we have observed in

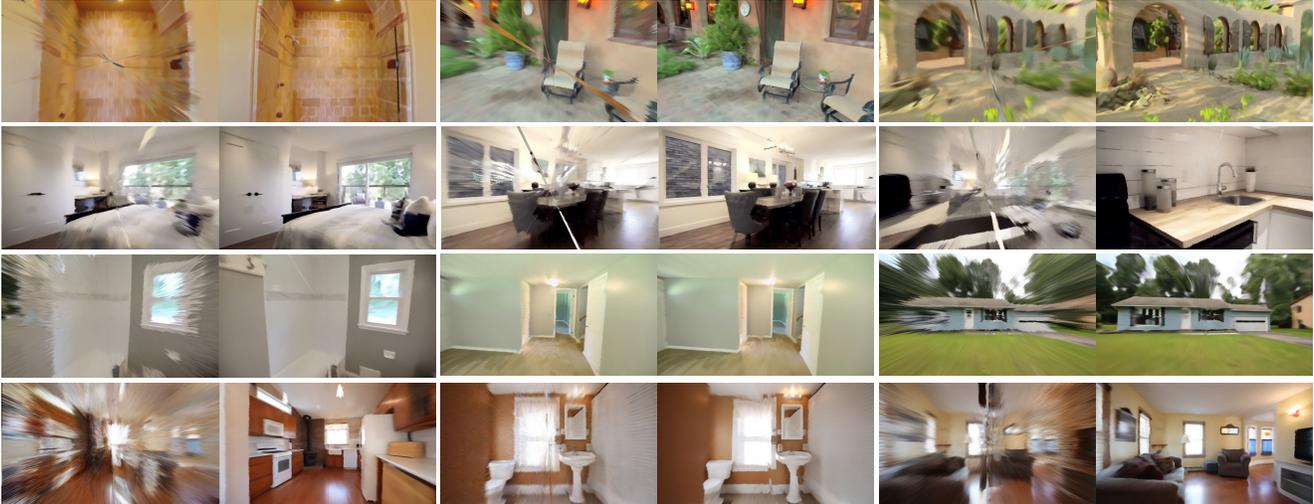


Figure 7. Qualitative comparison on Scenes from RealEstate10K dataset between mip360 (left image) and our method *Strata-NeRF* (right image) in a pair. Each row represents a scene in RealEstate10K and each pair represents a level in that scene. Our method outperforms and produce good quality novel views compared to mip360.

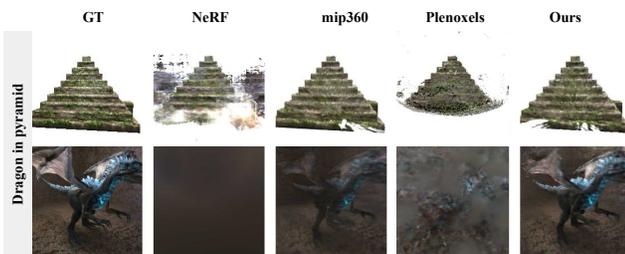


Figure 8. (From top to bottom) Qualitative results on the proposed synthetic datasets. Each row represents a novel view from each level of the stratified scene. The ground-truth view is shown in Column 1. Compared to prior works (Column 2-4) our method's (Column 5) renderings are more similar to the ground-truth.

Figure 9 that the baseline method fails on some of test images, hence we also compare the methods in worst care scenarios. The worst-case analysis describes a method's worst performance on the dataset. The worst case analysis is particularly useful to detect the shortcomings of the methods. We present analysis in two categories: (a) histogram distribution for each metric on the test set, and (b) qualitative comparison of the worst-case scenario for our method on PSNR metric.

Figure 9 compares PSNR histogram plots on test-set views for the “**Cube-Sphere-Monkey**” scene. We can see that the mip360 approach performs poorly on PSNR and ranks low on practically all stratification levels. This supports our argument that the mip360 approach produces artifacts in such stratified scenes. For our method, the PSNR distributions are on the right. This implies that the novel views on test-set from our method will not be having serious artifacts in most cases, demonstrating its reliability.

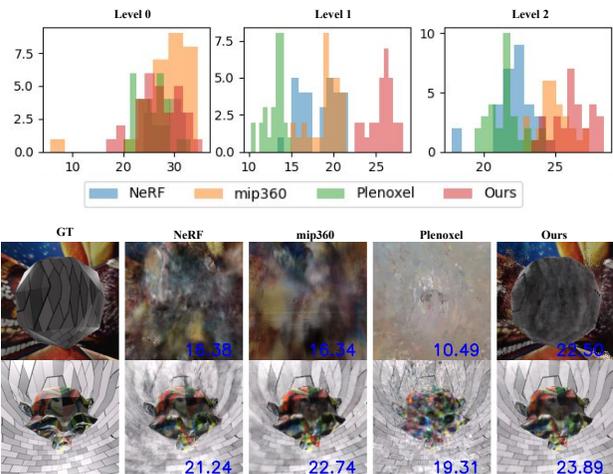


Figure 9. (Top Row) Comparison of histogram plots for the test-set for PSNR on “**Cube-Sphere-Monkey**”. Note how distribution of our our method is always towards the right compared to other methods.  $x$  – axis denote metric value and  $y$  – axis denotes the frequency. A qualitative comparison of our method's worst-case PSNR results. PSNR is present at the bottom of the result image.

Images in Figure 9 depict the qualitative results for the worst-case PSNR instances. All methods perform well in level 0. Hence, we are discussing interior levels which are level 1 and level 2. Other approaches fail in the worst-case scenario for our method at level 1. The outputs from NeRF, mip360 and Plenoxel are visually impaired. At level 2, our method has less blur compared to other approaches. These findings demonstrate that our method is better suited to represent stratified scenes than others.

**Ablation Studies.** To analyse our proposed method, we present an ablation on the size of the vector codebook in

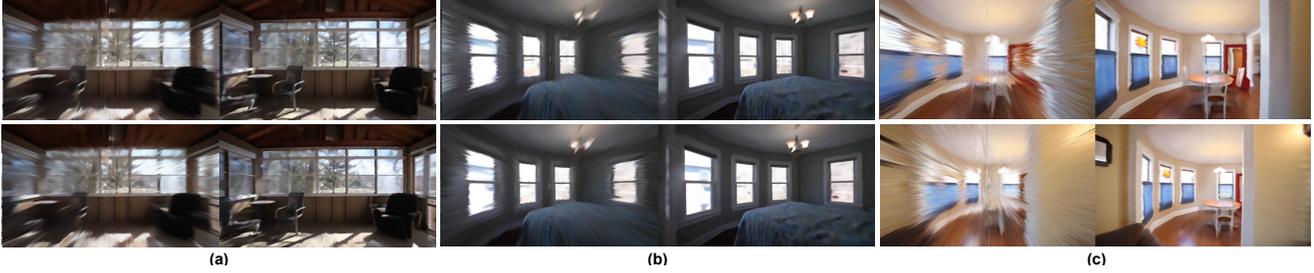


Figure 10. Novel-views from different levels of 'Real Estate Video Tour 7 Rutledge Ave Highland Mills NY 10930 Orange County NY' scene in Real Estate 10K dataset. The two rows are from two-different view-points.

Table 7. Quantitative comparison of our model and baseline on Synthetic Six Layer Scene.

	Metrics	Level 0	Level 1	Level 2	Level 3	Level 4	Level 5
mip360 [2]	PSNR $\uparrow$	22.215	16.183	15.084	12.012	21.813	21.539
	SSIM $\uparrow$	0.777	0.442	0.510	0.344	0.817	0.647
Ours	PSNR $\uparrow$	<b>23.889</b>	<b>21.449</b>	<b>21.456</b>	<b>24.095</b>	<b>28.283</b>	<b>21.898</b>
	SSIM $\uparrow$	<b>0.833</b>	<b>0.681</b>	<b>0.685</b>	<b>0.722</b>	<b>0.883</b>	<b>0.686</b>

Table 8. Quantitative results on "Cube-Sphere-Monkey" scene for ablation on size of the vector codebook in Latent Generator.

Size	PSNR $\uparrow$		SSIM $\uparrow$		LPIPS $\downarrow$	
	Level 0	Level 1	Level 0	Level 1	Level 0	Level 1
<b>512</b>	<b>29.5458</b>	26.3497	<b>0.8743</b>	0.7395	0.1675	<b>0.4899</b>
<b>1024</b>	29.4834	26.1715	0.8701	<b>0.7489</b>	<b>0.1367</b>	0.5163
<b>4096</b>	28.4609	<b>27.8274</b>	0.8628	0.7342	0.1776	0.5027

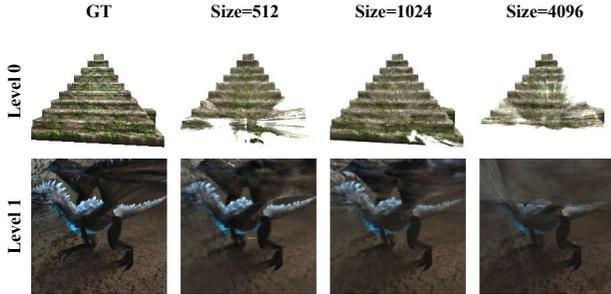


Figure 11. Comparisons of different codebook size on "Dragon in Pyramid" scene for different vector-codebook sizes. Note at size=1024 we achieve the best results with less artifacts.

our latent generator. Table 8 shows the ablation for the size of the vector codebook on the "Coffee Shop" dataset. We trialed with codebook sizes of 512, 1024 and 4096. We found that size 1024 provides optimal performance. As shown in Figure 11, increasing the codebook size induces haziness in the generated novel views, while decreasing the size creates white artifacts in level 0. As a result, we fix the size 1024 for all of our synthetic experiments. Whereas for RealEstate10K dataset we find that codebook size of 4096 produces the optimal tradeoff of results across levels, as it contains more number of levels and details. We further discuss the key architectural design choices for Latent Generator and Latent Router modules in Appendix E.5.

**No. of levels:** To further test the efficacy of our method on higher number of levels, we created a "Simple Geometry" scene consisting of primitive geometry shapes like cube and spheres. More details are in the supplementary material. Table 7 displays the results for both the baseline and our approach across a six levels stratified scene. The average PSNR/SSIM for the mip360 baseline is **15.35 / 0.487**, while our method achieved PSNR/SSIM of **23.54 / 0.754**

which improves PSNR and SSIM by **53.35 %** and **54.83 %** respectively. This shows that our method performs better on increasing number of levels when compared with the baseline method. These observations also hold true for scenes in the RealEstate10K dataset as shown in Table 5.

## 7. Conclusion

In this work, we focus on the problem of modelling the 3D representation of a stratified and hierarchical scene, implicitly through a single neural field. For this, we propose *Strata-NeRF*, which models scenes with stratified structures by introducing a VQ-VAE-based latent generator to implicitly learn the distribution of latent space of input 3D locations and condition the neural radiance field with the latent code generated from this distribution. We also introduce a new synthetic dataset with stratified-level scenes and use it to analyse various existing approaches. Through quantitative, qualitative, and worst-case analysis on this dataset, we show that *Strata-NeRF* has a more stable 3D representation than the other methods. Further, the improvements due to *Strata-NeRF* also generalize to real-world RealEstate10K dataset, where it outperforms baselines by a significant margin establishing a new state-of-the-art. We believe designing a new volume rendering equation for modelling complex stratified scenes is a good direction for future work.

**Acknowledgement.** This work was supported by Samsung R&D Institute India, Bangalore, PMRF and Kotak IISc AI-ML Centre (KIAC). Srinath Sridhar was partly supported by NSF grant CNS-2038897

# Appendix

## Table of Contents

---

<b>A Introduction</b>	<b>10</b>
<b>B Synthetic Dataset Details</b>	<b>10</b>
B.1. Cube-Sphere-Monkey . . . . .	10
B.2. Coffee Shop . . . . .	10
B.3. Bhutanese House . . . . .	10
B.4. Dragon In Pyramid . . . . .	10
B.5. Buddhist Temple . . . . .	10
<b>C Real Dataset</b>	<b>11</b>
<b>D Implementation Details</b>	<b>11</b>
D.1. Choice of Training Configuration File	11
<b>E Additional Experiments</b>	<b>12</b>
E.1. RealEstate10K [74] scene - Motivation Experiment . . . . .	12
E.2. Comparison with InstantNGP [39] and TensoRF [8] . . . . .	13
E.3. Comparison with level-wise radiance fields. . . . .	13
E.4. Ablation on Vector-Codebook Size . . . . .	13
E.5. Architectural Design Choices. . . . .	13
E.6. Why shared codebooks are important? . . . . .	14
E.7. Experiments on the standard novel-view synthesis dataset. . . . .	14
E.8. Number of Views . . . . .	15
E.9. Out of Distribution Views . . . . .	15
E.10. Additional Results . . . . .	15
E.11. Impact of Image-Resolution on training. . . . .	15

---

### A. Introduction

We present additional results and other details related to our proposed method : Strata-NeRF. We elaborate on the proposed synthetic stratified dataset in Appendix B. We give the implementation details in Appendix D. Then, we present additional ablation study and results in Appendix E.

### B. Synthetic Dataset Details

Figure 12 shows the representation of each level of each scene. Table 9 shows the level-wise split for each scene.

#### B.1. Cube-Sphere-Monkey

This dataset consists of simple geometric entities such as a cube, sphere and a monkey mesh provided in Blender [10]. Figure 12 illustrates the layout of this scene. *Cube* is at level 0, *Sphere* is at level 1 and *Monkey* is at the innermost level. The texture for *Cube* is an image generated from Stable Diffusion demo [18]. We sample camera poses from the curved surface of a hemisphere for the outer cube and from the curved surface of a sphere for the inner levels.

#### B.2. Coffee Shop

This dataset mimics an actual coffee shop setup inside another shopping complex. The outermost level consists of concrete walls. At level 1, i.e. when one enters the shopping complex, there is regular flooring and a concrete ceiling. Here, we also notice the exterior walls of our coffee shop. At level 2; i.e., inside the coffee shop; there is a layout with a counter, menu board and a table for visitors. All these scenes are composited with the help of Blender [10]. We sample camera poses from the curved surface of a hemisphere for all the levels.

#### B.3. Bhutanese House

A typical household setting inspired us to create this dataset. A typical residence features a table in the living room. In most cases, a decorative object is kept on the table. For the structure of the house, we choose a Bhutanese house model. The exterior of this structure is level 0. At level 1, i.e., inside the house, there are chairs, tables and other household items in the living room. At level 2, we have a glass bottle with a ship. We sample camera poses from the curved surface of a hemisphere. For level 2, we capture around the glass bottle on the circular table.

#### B.4. Dragon In Pyramid

This dataset captures a fantastical world filled with pyramids and dragons. We use a model of a *Mayan pyramid* as the outer structure. Inside the pyramid, we place a model of a dragon. Thus, this scene has two levels: 1.) the outer walls of the *Mayan pyramid* and 2.) the dragon residing inside the pyramid. All the camera poses are sampled from the curved surface of different hemispheres.

#### B.5. Buddhist Temple

This scene depicts an archaeological site or a typical monument location. We select a Buddhist temple to represent this scene. Two levels indicate the nearby rooms inside the structure in this context. Level 0 represents the outer structure of the monument, Levels 1 contains a bronze statue in the center of the monument, and Level 2 contains a Buddha statue mounted to the wall of one room.

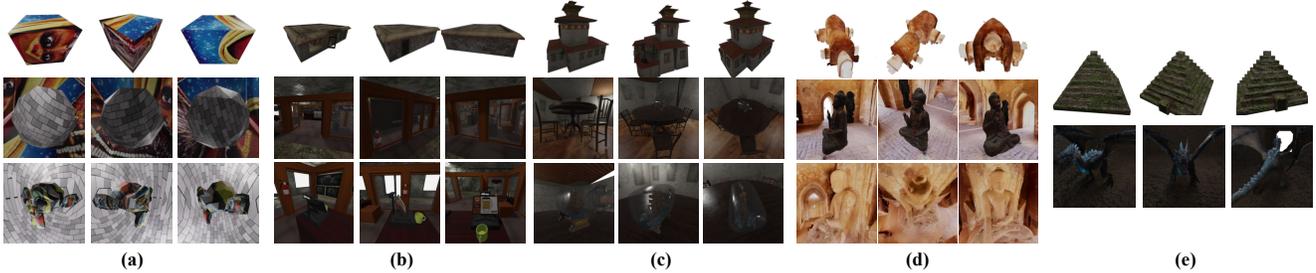


Figure 12. (a) Cube-Sphere-Monkey, (b) Coffee Shop, (c) Bhutanese House, (d) Buddhist Temple and (e) Dragon In Pyramid. Representative images for each level.

Table 9. **train-val-test** level-wise split for each scene.

Scene	Split	Level 0	Level 1	Level 2
Cube-Sphere-Monkey	train	30	30	30
	val	30	30	30
	test	30	30	30
Coffee Shop	train	30	30	30
	val	15	15	15
	test	15	15	15
Bhutanese House	train	30	30	30
	val	15	15	15
	test	15	15	15
Buddhist Temple	train	30	20	20
	val	15	10	10
	test	15	10	10
Dragon In Pyramid	train	30	30	-
	val	15	15	-
	test	15	15	-

## C. Real Dataset

We evaluate our method on real-world scenes as well. We choose RealEstate10K [74] dataset, which contains camera poses corresponding to camera frames from video-clips extracted from Youtube videos. The camera poses are obtained by running SLAM and bundle adjustment algorithm over these large videos. To create a “stratified” scene from this dataset, first we cluster video clips belonging to same Youtube video using the video token provided in the ground-truth files. Then we extracted camera frames and pose as per the timestamp information provided in the ground-truth files. The extracted camera pose for each video clip from a scene were already aligned with respect to a common coordinate system. We removed the video clips which had any dynamic motion within them. We extracted four scenes which are “Spanish Colonial Retreat in Scottsdale Arizona” [47], “139 Barton Avenue Toronto Ontario” [59], “31 Brian Dr Rochester NY” [4] and “7 Rutledge Ave Highland Mills” [24].

## D. Implementation Details

**Architecture Details.** We provide architectural details of the “Latent Generator” and “Latent Router” networks in Figure 13 and 14 respectively.

**Training.** We use Adam [29] optimizer with hyperparameters  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ,  $\epsilon = 1e^{-6}$  and initial learning rate = 0.002. Further, the learning rate is log-linearly interpolated such that learning rate = 0.00002 at maximum steps. Additionally, there are 512 warmup steps. Distortion loss proposed in Mip-NeRF 360 [2] is switched off for the blender datasets as proposed by the authors. We use one proposal MLP and one NeRF MLP. We weight the loss for “Latent Generator” with value  $\lambda_2 = 0.1$ .

**Implementation.** Our implementation is based on Mip-NeRF 360 [2] which uses JAX [5] framework.

### D.1. Choice of Training Configuration File

The dataset described in Section B is created using Blender [10]. This dataset has white background for the level 0. Barron *et al.* [2] uses “blender\_256.gin” file for the blender scenes proposed in NeRF [38] which are small in size compared to our scenes. This configuration file does not work for the scenes we proposed in Appendix B. Hence,

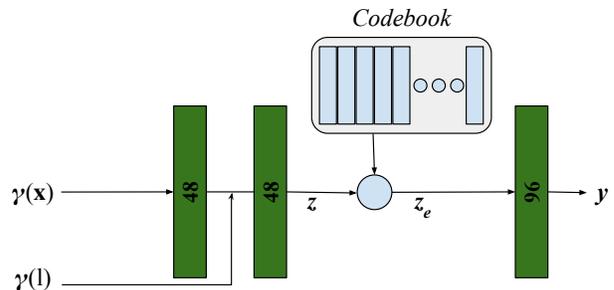


Figure 13. A diagram of “Latent Generator” network. This network takes position-encoded 3D point  $\gamma(x)$  and position-encoded camera level  $\gamma(l)$ . This is passed through the encoder block to get  $z$  which is then matched to the nearest latent in the codebook to get  $z_e$ .  $z_e$  is passed through decoder block to reconstruct the position-encoded 3D point  $y$ .

Table 10. Performance on the *Dragon In Pyramid* dataset between two configuration files. We observe that “360.gin” works much better than the other configuration file.

Config	Level 0			Level 1			Total		
	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS
Blender	5.5654	0.3717	0.6252	22.9489	0.6320	0.5844	14.2571	0.5018	0.6048
360	<b>30.8758</b>	<b>0.9006</b>	<b>0.1367</b>	<b>24.3890</b>	<b>0.7054</b>	<b>0.5163</b>	<b>27.6324</b>	<b>0.8030</b>	<b>0.3265</b>

we use “360.gin” and alter the dataset type field in the configuration file.

Table 10 shows the quantitative comparison of the above mentioned configuration files on *Dragon In Pyramid* dataset. We observe that the “360.gin” configuration beats the “blender\_256.gin” in all the levels. Figure 15 compares the qualitative results of these two configuration files. We notice that the novel views from “blender\_256.gin” are inferior in quality compared to “360.gin” configuration. “360.gin” configuration has better performance because of the contract function proposed by Barron [2]. The contract function is defined as follows:

$$\text{contract}(x) = \begin{cases} x, & \|x\| \leq 1 \\ (2 - \frac{1}{\|x\|}) \left( \frac{x}{\|x\|} \right), & \text{otherwise} \end{cases} \quad (5)$$

This contract function maps input coordinates onto a ball of radius 2. Effectively, a large range is bounded inside a radius of  $2m$ . This is the reason why “360.gin” configuration is better for large blender scenes. Hence, we use this configuration file for all the scenes other than “Cube-Sphere-Monkey”.

## E. Additional Experiments

### E.1. RealEstate10K [74] scene - Motivation Experiment

We presented motivation of our work on a synthetic scene “Dragon In Pyramid” in Section 4 in the main paper. We observed that no artifacts are observed if individual mipNeRF-360 is trained for each level (level-wise) separately. We performed a similar experiment on the RealEstate10k [74] scene and observed artifact-free novel

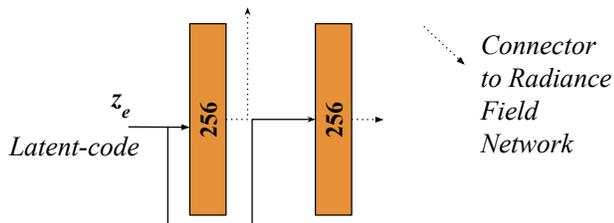


Figure 14. A diagram of “Latent Router” network. This network takes latent code  $z_e$  generated by the “Latent Generator” and connects it to the radiance field network after passing through linear layers.

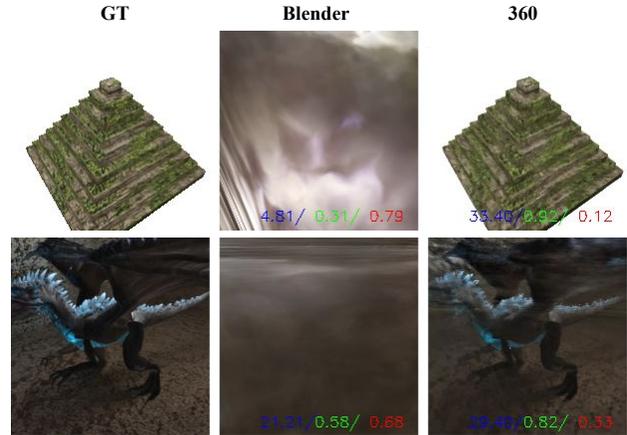


Figure 15. Qualitative comparison for different configuration files on *Dragon In Pyramid* scene. We observe that 360.gin configuration generates better results. Metrics PSNR, SSIM and LPIPS are color-coded at the bottom of the result image

Table 11. No. of training parameters (in millions) for level-wise mip360 and our method with two different codebook sizes 1024 and 4096 for different number of levels.

Levels	Level-Wise mip360	Ours (1024 codebook)	Ours (4096 codebook)
1	0.835	0.924	1.071
3	2.506	0.924	1.071
4	3.341	0.924	1.071
5	4.176	0.924	1.071
6	5.011	0.924	1.071

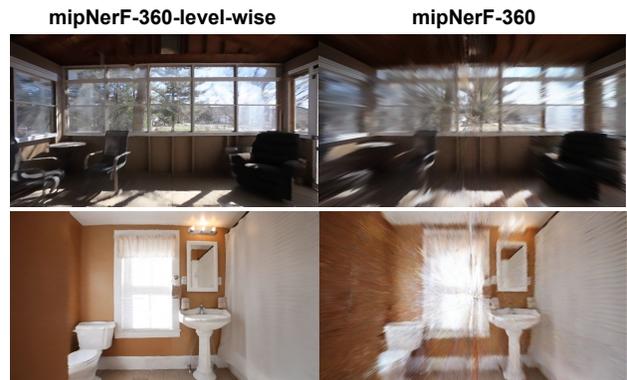


Figure 16. Analysis on “7 Rutledge Ave” scene from RealEstate10K [74] dataset. We present visual results from two levels. Note how artifacts appear in results from mipNeRF-360 (all levels are trained jointly) whereas when mipNeRF-360 is used for each level separately (level-wise) we observe no artifacts.

views from level-wise mipNeRF-360. Similar to the observation for synthetic scenes, if all levels are trained combinedly we observe the artifacts in the rendered novel-views as shown in Fig 16. Further, PSNR values in Tab. 12 for

Table 12. A quantitative comparison of mip360 (level-wise) and mipNeRF-360 (all views) on “7 Rutledge Ave”

Methods	Lv 0	Lv 1	Lv 2	Lv 3	Lv 4	Lv 5	Lv 6	Total
mipNeRF-360 (x7)	<b>24.20</b>	22.42	<b>26.72</b>	24.78	22.73	<b>27.41</b>	24.78	<u>24.25</u>
mipNeRF-360	19.53	18.33	23.52	17.00	18.82	19.73	21.60	19.62

Table 13. A quantitative comparison of InstantNGP [39] and TensorRF [8] on “7 Rutledge Ave”

Methods	Lv 0	Lv 1	Lv 2	Lv 3	Lv 4	Lv 5	Lv 6	Total
Instant-NGP	19.02	18.24	21.32	19.43	18.77	18.98	21.33	19.47
TensorRF	18.03	21.29	21.23	20.23	20.36	18.57	22.69	20.70
Ours	<b>22.84</b>	<b>25.14</b>	<b>24.83</b>	<b>25.67</b>	<b>25.15</b>	<b>23.10</b>	<b>26.75</b>	<b>25.04</b>

level-wise mipNeRF-360, with 7 radiance fields (x7) are higher compared to a single mipNeRF-360 for all-levels. This further substantiates our claim that a single mipNeRF-360 network is not able to learn all the stratified levels.

## E.2. Comparison with InstantNGP [39] and TensorRF [8]

**Synthetic Scenes.** We present qualitative comparison with InstantNGP [39] and TensorRF [8] in Fig. 17 and ???. These methods work well in the outermost level. But suffer from artifacts because of the stratified scenes in the inner levels. We observe this pattern consistently across all the synthetic scenes.

**RealEstate10K [74] dataset** Fig. 18 shows qualitative comparison on “7 Rutledge Ave” scene from RealEstate10K [74]. Our method generates novel-view without any artifact, whereas other methods have visible artifacts in the generated novel-views. Tab. 13 shows PSNR of the generated novel-views. Our method clearly outperforms InstantNGP [39] and TensorRF [8].

## E.3. Comparison with level-wise radiance fields.

One trivial solution for the proposed stratified setting is training mip360 individually for multi-view images in each level. We show that with increase in no. of levels, no. of training parameters increases linearly. Consider a mip360 network with width 256 and depth 8. We present variation of no. of training parameters in Table 11 for different number of levels. Our method’s training parameter requirement doesnot increase linearly as it does in level-wise mip360.

For comparison, on “Spanish Colonial Retreat” scene, mipNerf-360 takes *5h 30m* to train, while our method, with a vector-codebook size of 1024, takes *6h 20m* for 150k iterations on a single NVIDIA RTX 3090 GPU.

## E.4. Ablation on Vector-Codebook Size

We present more results on *Coffee Shop*, *Bhutanese House* and *Buddhist Temple* for the ablation : *Size of the vector-codebook in “Latent Generator”*. We tried with three sizes : 512, 1024 and 4096. Table 14 and 15 shows the quantitative results for the mentioned datasets. We observe

Table 14. Performance on the *Coffee Shop* dataset for different sizes of the vector codebook. **Best** results are marked in bold and **Second-best** results are underlined.

Size	Level 0			Level 1			Level 2			Total		
	PSNR	SSIM	LPIPS									
512	24.4768	0.8605	0.2049	28.0758	0.8257	0.3632	<b>33.7944</b>	<u>0.9306</u>	<b>0.2003</b>	28.7824	0.8723	0.2561
1024	<b>26.4497</b>	<b>0.8803</b>	<b>0.1936</b>	<b>28.6387</b>	<b>0.8403</b>	<b>0.3449</b>	33.2695	0.9254	0.2243	<b>29.4526</b>	<b>0.8820</b>	<u>0.2543</u>
4096	<u>25.3534</u>	<u>0.8729</u>	<u>0.1995</u>	28.4341	0.8383	0.3539	<u>33.6062</u>	<b>0.9316</b>	<u>0.2025</u>	<u>29.1312</u>	0.8809	<b>0.2520</b>

Table 15. Performance on the *Buddhist Temple* dataset for different sizes of the vector codebook. **Best** results are marked in bold and **Second-best** results are underlined.

Size	Level 0			Level 1			Level 2			Total		
	PSNR	SSIM	LPIPS									
512	<u>27.3121</u>	<u>0.8881</u>	<u>0.1861</u>	25.3407	0.7619	0.362	<u>25.4983</u>	<u>0.7476</u>	<u>0.3691</u>	<u>26.2306</u>	<u>0.8119</u>	<u>0.2886</u>
1024	<b>27.5529</b>	<b>0.8935</b>	<b>0.1775</b>	<b>27.3453</b>	<b>0.7894</b>	<b>0.3240</b>	<b>25.5956</b>	<b>0.7717</b>	<b>0.3456</b>	<b>26.9343</b>	<b>0.8289</b>	<b>0.2674</b>
4096	20.9017	0.8075	0.2680	<u>27.0011</u>	<u>0.7856</u>	<u>0.3340</u>	23.4656	0.7189	0.3853	23.3769	0.7759	0.3204

Table 16. Ablation studies on the key design choices for the proposed method. **D1**: Disable second router in LR, **D2**: Disable first router in LR, **D3**: Remove LR and directly concatenate generated embedding with the positional encoding and **D4**: Replace VQ-VAE with VAE in LG. Acronyms D1, D2, D3, D4 are explained in more detail in Appendix E.5

	D1	D2	D3	D4	Ours
<b>Synthetic</b>	26.04	27.34	27.41	26.96	<b>28.25</b>
<b>RealEstate10K</b>	23.79	24.24	23.79	20.99	<b>24.75</b>

that vector codebook of size 1024. gives us the overall best results.

## E.5. Architectural Design Choices.

The proposed method consists of Latent Generator (LG) and Latent Router(LR) as shown in Figure 4 in the main paper. Latent Generator(LG) and Latent Router(LR) are described in Section 5.1 and 5.2 respectively in the main paper. To further motivate this choice of the architecture, we discuss the following design choices for the proposed method:

1. Disabling the second router in **LR: D1**
2. Disabling the first router in **LR: D2**
3. removing **LR** and directly concatenating the generated embedding to the input positional encoding : **D3**
4. Replacing the VQ-VAE block with the VAE block in **LG : D4**

We present overall results for synthetic and RealEstate10K scenes in Tab. 16. We conclude that using two parallel dense layers is better than an individual dense layer in **LR**. Further, we observe that how using Latent Router is better than directly concatenating the generated embedding with the input positional embedding. Similarly, the VAE version of our method underperforms the discrete VQ-VAE used in our method.

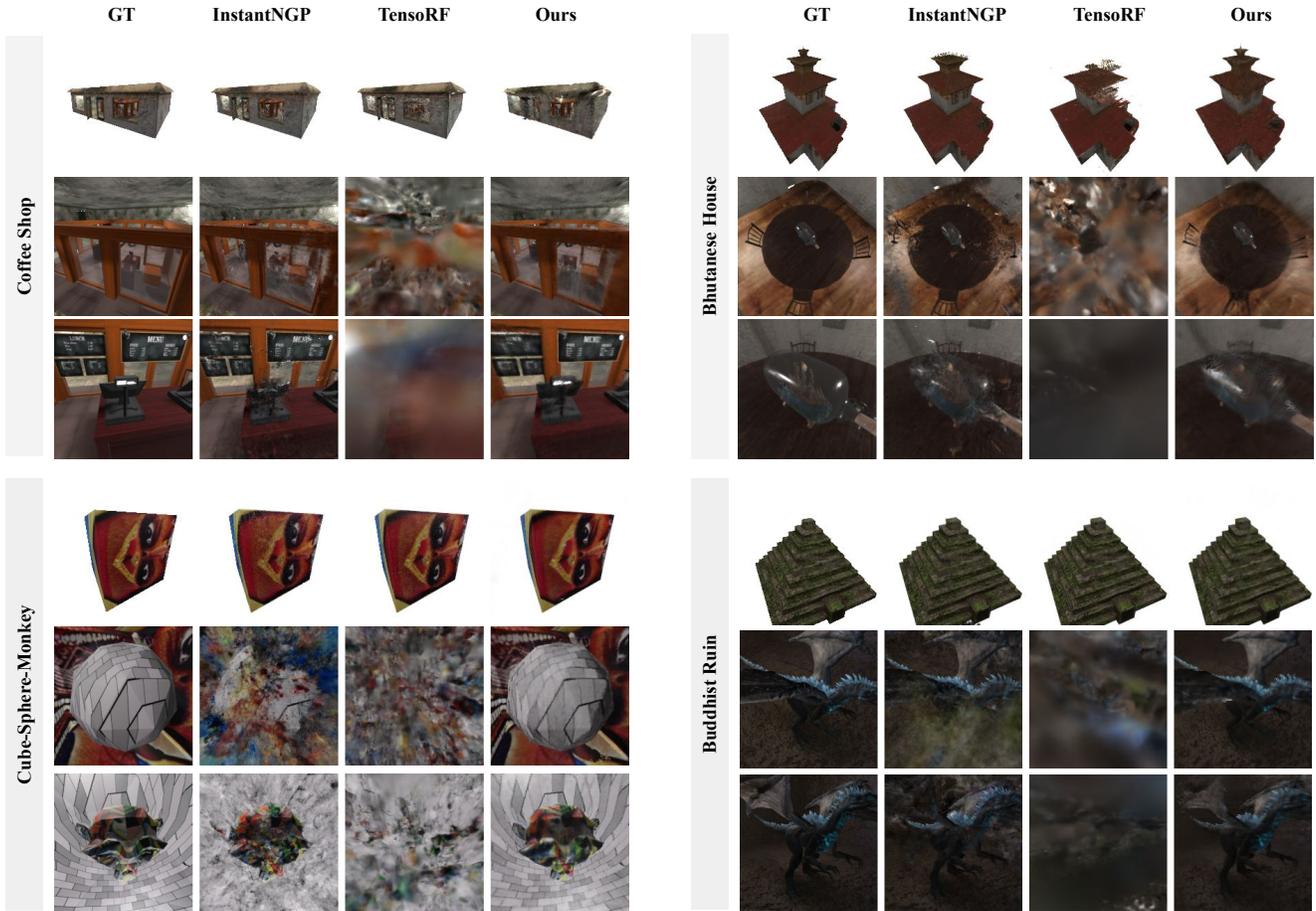


Figure 17. Qualitative Comparison on synthetic dataset for InstantNGP [39] and TensoRF [8]



Figure 18. Qualitative Comparison on “7 Rutledge Ave” scene from RealEstate10K [74] dataset. The novel-view generated from our method is better than InstantNGP [39], TensoRF [8] and mipNeRF-360 [2]

Table 17. Quantitative Comparison on “7 Rutledge Ave”

	21.03	23.54	24.15	23.85	22.83	22.64	25.41	23.53
<b>Ours-Ind.</b>	21.03	23.54	24.15	23.85	22.83	22.64	25.41	23.53
<b>Ours</b>	<u>22.84</u>	<u>25.14</u>	<u>24.83</u>	<u>25.67</u>	<u>25.15</u>	<u>23.10</u>	<u>26.75</u>	<u>25.04</u>

### E.6. Why shared codebooks are important?

We provide another ablation by creating independent code-book vectors for different levels : “Ours-Ind.”. In our method, codebooks are shared between level which yield better results. This is natural as walls, etc. are shared between levels in the scene.

### E.7. Experiments on the standard novel-view synthesis dataset.

We train the “garden” scene from the mipNeRF-360 dataset by treating it as a single-level scene. We achieved a PSNR of 26.40 on the test dataset, while mipNeRF-360 reports a PSNR of 26.98. We achieve an average PSNR of 33.21 across all NeRF-synthetic scenes, while mipNeRF-360 achieves 33.09. Our proposed method performs comparably on these datasets, despite being designed for stratified scenes.

Table 18. Performance on the *Dragon In Pyramid* dataset for different number of views in the dataset. **Best** results are marked in bold.

		Level 0			Level 1			Total		
		PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS	PSNR	SSIM	LPIPS
1x Views	mip360	<b>30.8758</b>	<b>0.9006</b>	<b>0.1367</b>	24.3890	0.7054	0.5163	27.6324	0.8030	0.3265
	Ours	29.4773	0.8700	0.1699	<b>26.1722</b>	<b>0.7489</b>	<b>0.4573</b>	<b>27.8248</b>	<b>0.8095</b>	<b>0.3136</b>
2x Views	mip360	<b>29.5127</b>	<b>0.8436</b>	<b>0.1830</b>	26.2172	0.7245	0.4627	27.8650	0.7841	0.3228
	Ours	29.1104	0.8099	0.2176	<b>27.4282</b>	<b>0.7661</b>	<b>0.4244</b>	<b>28.2693</b>	<b>0.7880</b>	<b>0.3210</b>
3x Views	mip360	<b>31.1511</b>	<b>0.8764</b>	<b>0.1715</b>	26.5231	0.7239	0.4638	28.8371	0.8001	0.3176
	Ours	30.5436	0.8461	0.1882	<b>27.4354</b>	<b>0.7693</b>	<b>0.4385</b>	<b>28.9895</b>	<b>0.8077</b>	<b>0.3134</b>

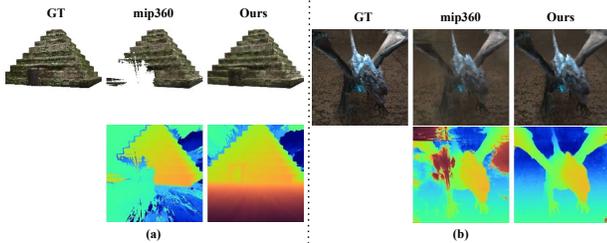


Figure 19. Qualitative Results for  $2\times$  views on *Dragon In Pyramid* scene. Observe that our results have less artefacts and much smoother depth maps.

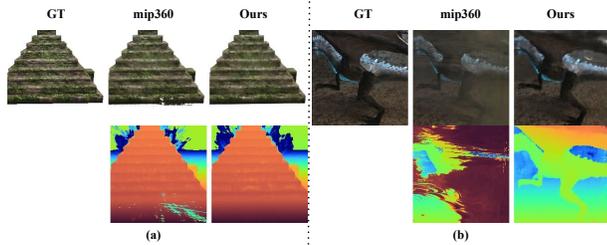


Figure 20. Qualitative Results for  $3\times$  views on *Dragon In Pyramid* scene. Observe that our results have less artefacts.

## E.8. Number of Views

We present here another ablation which evaluates the effect of increasing number of views for a scene. Table 18 shows quantitative results on *Dragon In Pyramid* scene by increasing number of views  $2\times$  and  $3\times$ . Note that  $2\times$  views mean that train, validation and test views will be doubled. We observe that as number of views are increased, overall metrics improves in both mip360 [2] and our method. Further, we compare qualitative performance of our method with mip360 [2] with increased number of views in Figure 19 and 20. We observe that quality of depth map is much better in our method. Also, generated novel views from our method has less artefacts.

## E.9. Out of Distribution Views

The training set’s views are uniformly sampled from the curved surface of a hemisphere with the camera’s  $z$  - axis always pointing towards the subject. Out-of-distribution (OOD) is any new view that does not lie on this hemisphere and whose  $z$  - axis is not necessarily aligned with the subject. We investigated the quality of novel view synthesis for

OOD views. We apply a random rotation and translation to the camera pose in the test set to produce OOD camera poses. A random translation value is sampled uniformly between  $(10cm, 10cm)$ , which is then used to translate the camera position along its  $z$  - axis. We randomly choose the rotation axis and angle from  $(-45^\circ, 45^\circ)$  for random rotation and change the current pose with this transformation. Figure 21 shows the novel views and their corresponding depth maps. The depth map shows that our technique regularises the 3D geometry significantly better than other methods. Furthermore, the depth map quality is substantially better, which aids our method in producing non-blurry results.

## E.10. Additional Results

We provide more results for the Out Of Distribution views in Figure 22. Further, we provide a sequence of generated novel views for *Cube-Sphere-Monkey* in Figure 23 and a sequence of depth maps for the *Buddhist Temple* in Figure 24. There are distinct artefacts in column one and three in Figure 22(a), column one in 22(b) and column three in 22(b). We compare the generated depth maps in Figure 22 and Figure 24. We observe that the depth maps from our method are smooth and have less artefacts than Mip-NeRF 360 [2]. Notice the collapse in floor of the *Buddhist Temple* scene in Figure 24. From these results, it’s clear that the generated novel views from our method has less artefacts and better 3D representation of such stratified scenes.

## E.11. Impact of Image-Resolution on training.

On  $800 \times 800$  resolution for “Cube-Sphere-Monkey” scene, mipNeRF-360 achieves an overall PSNR of 23.17 and our method achieves 26.41. This is similar to behavior observed on low-resolution ( $200 \times 200$ ) and high-resolution RealEstate10K.

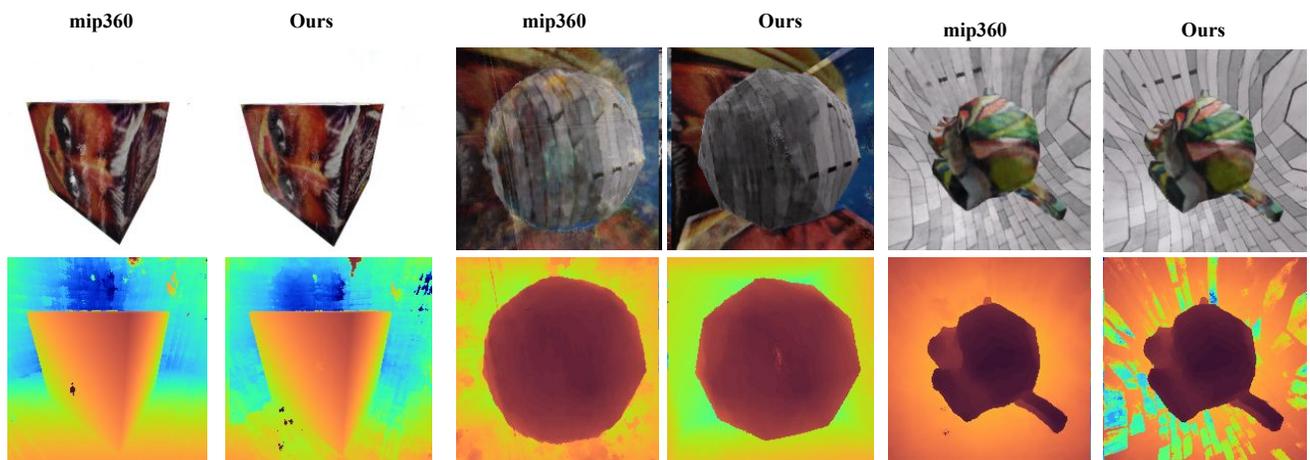
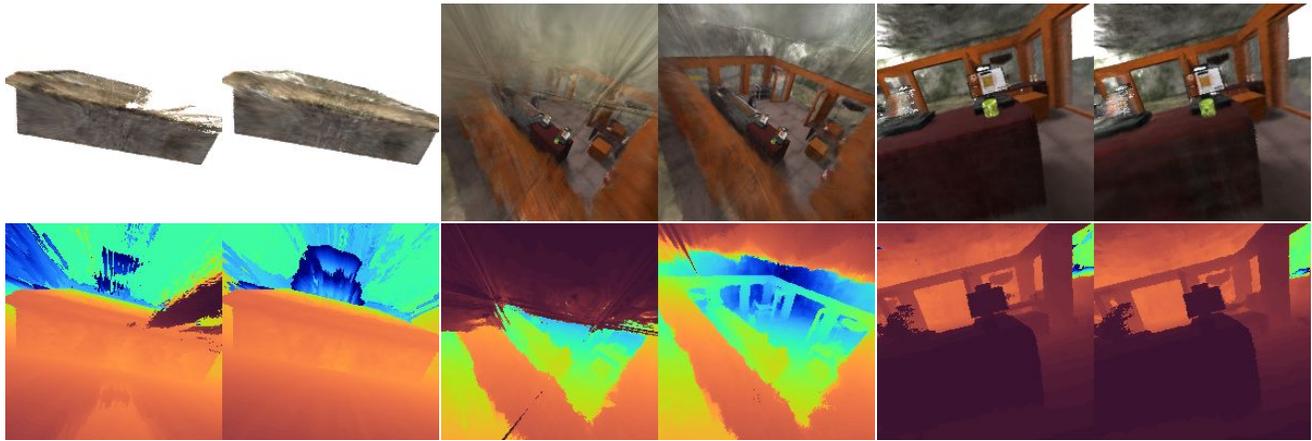
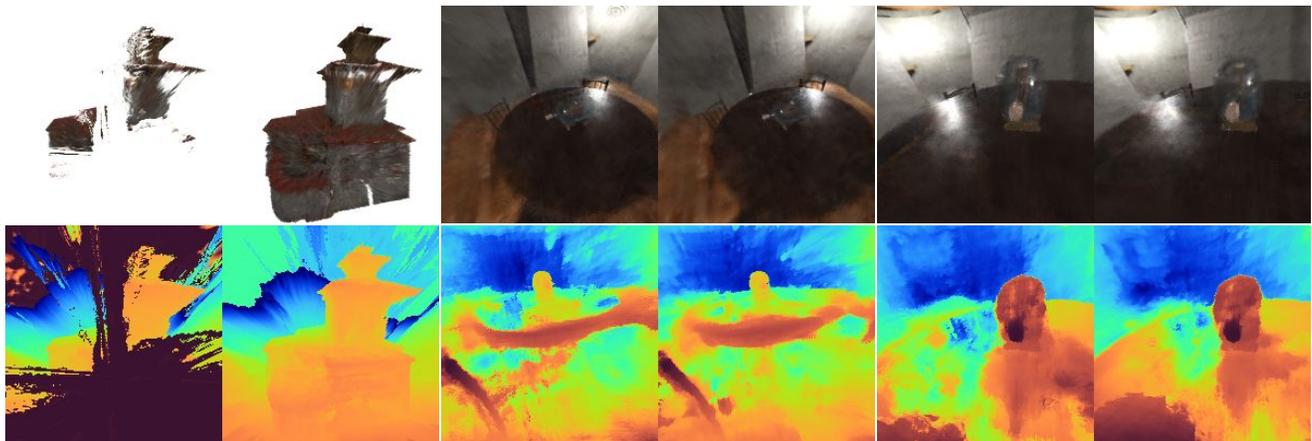


Figure 21. Qualitative comparison on OOD views. (**Top Row**) Generated novel views. (**Bottom Row**) Corresponding depth map. Check the quality of depth maps in inner levels for our method.



(a)



(b)



(c)

Figure 22. Out of distribution views for (a) Coffee Shop, (b) Bhutanese House and (c) Dragon In Pyramid Scene. **Odd** columns are results from Mip-NeRF 360 [2] and **even** columns are results from our method. We observe that generated novel views from our method has less artefacts and better depth maps. Check the clarity in claws of dragon in last column of (c).

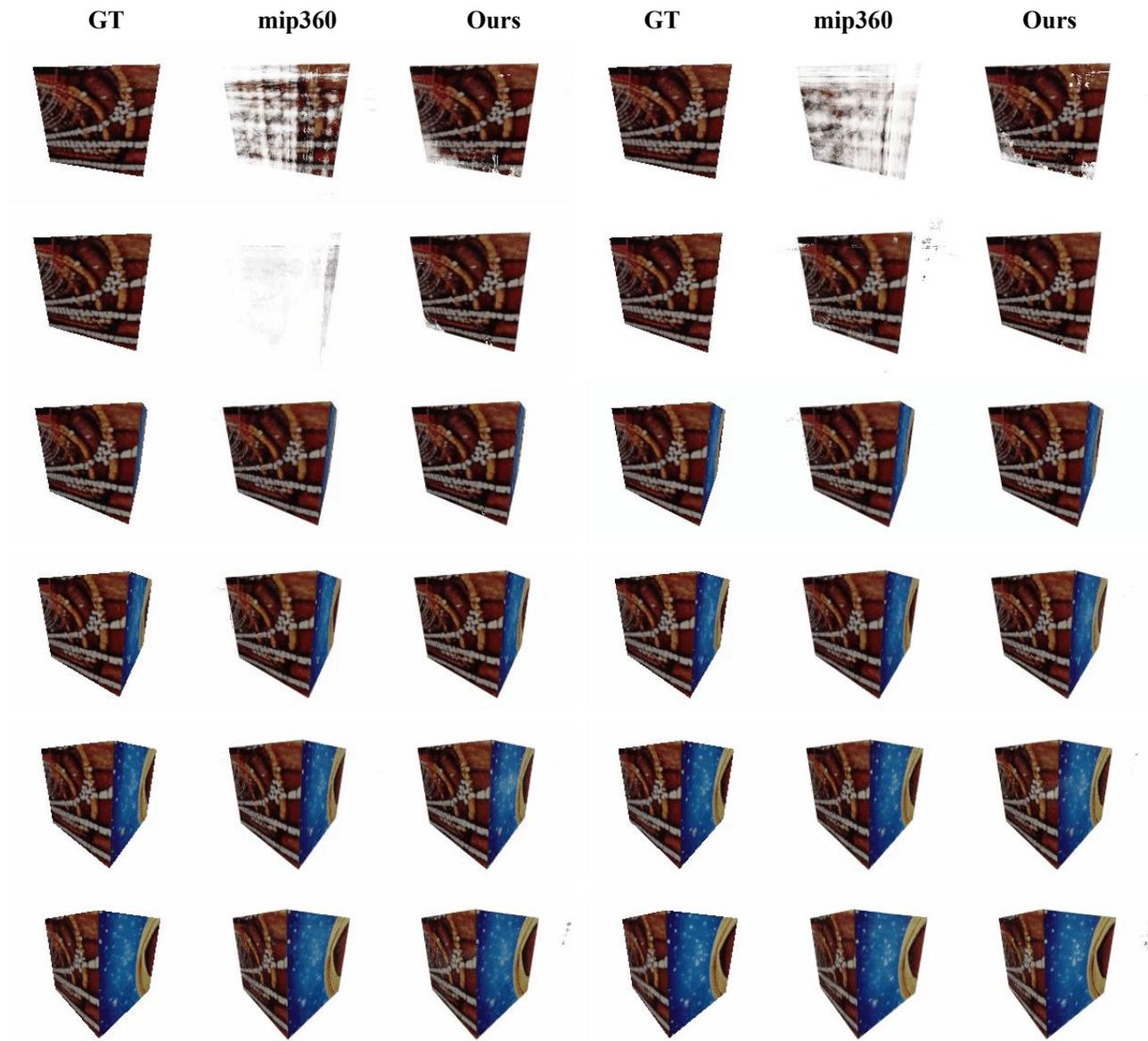


Figure 23. Sequence of generated novel views for Level 0 of *Cube-Sphere-Monkey* scene. Please note that sequence is represented in zig-zag pattern. The generated novel views from our method has less artefacts. **Please check the video provided in the supplementary material to appreciate our results better.**

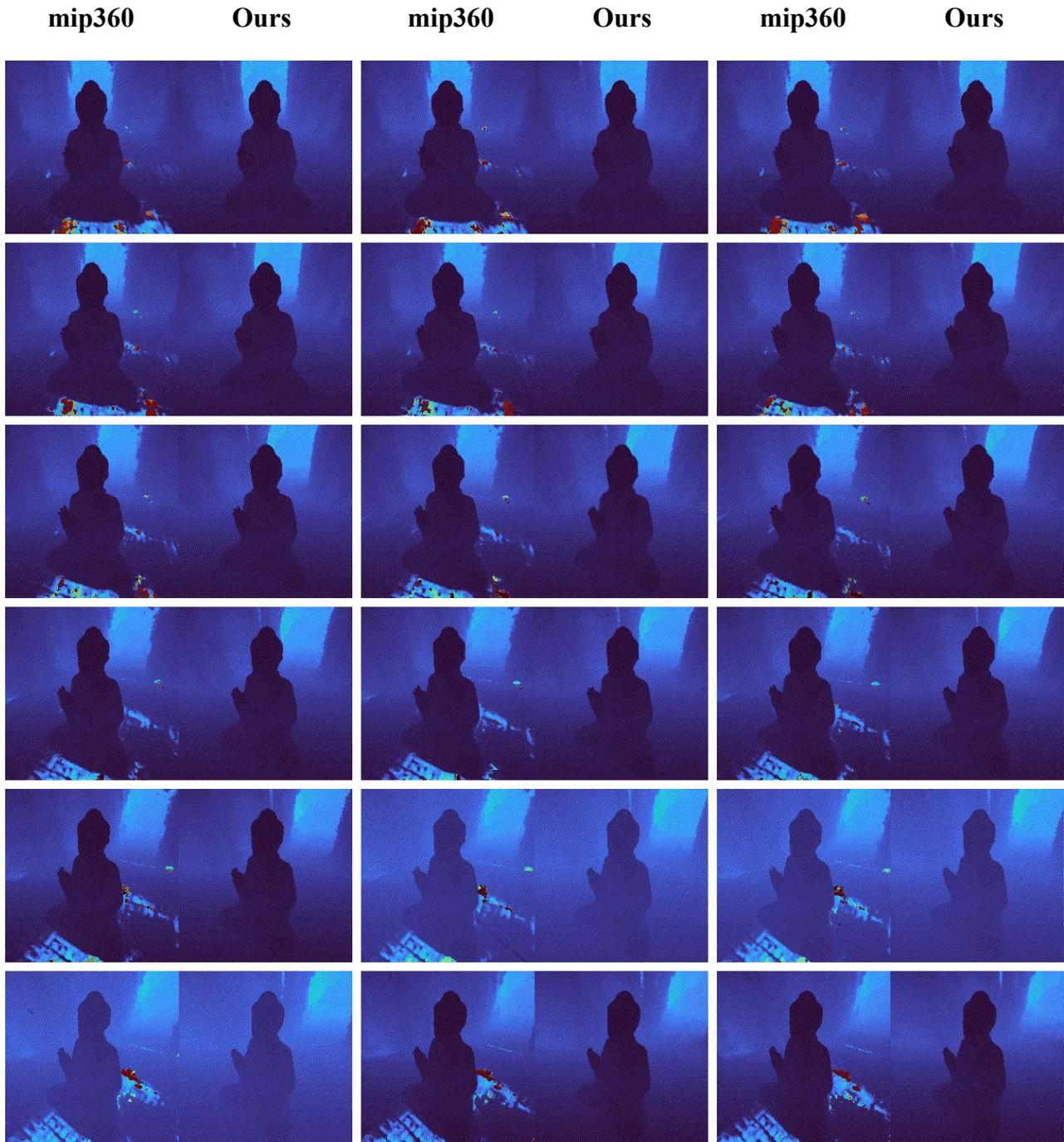


Figure 24. Sequence of depth maps of generated novel views for Level 1 of *Buddhist Temple* scene. Please note that sequence is represented in zig-zag pattern. We observe that there is a collapse in the floor region for output from mip360 [2] output. Whereas, our method generates smooth depth maps. **Please check the video provided in the supplementary material to appreciate our results better.**

## References

- [1] Jonathan T Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. [2](#), [3](#)
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [9](#), [11](#), [12](#), [14](#), [15](#), [17](#), [19](#)
- [3] Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition. *arXiv preprint arXiv:2008.03824*, 2020. [3](#)
- [4] birdhousemediatv. 139 barton avenue, toronto, ontario. [11](#)
- [5] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, et al. Jax: composable transformations of python+ numpy programs. *Version 0.2*, 5:14–24, 2018. [11](#)
- [6] Shengqu Cai, Anton Obukhov, Dengxin Dai, and Luc Van Gool. Pix2nerf: Unsupervised conditional p-gan for single image to neural radiance fields translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [3](#)
- [7] Eric R Chan, Marco Monteiro, Petr Kellnhofer, Jiajun Wu, and Gordon Wetzstein. pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [3](#)
- [8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European conference on computer vision (ECCV)*. Springer, 2022. [5](#), [6](#), [7](#), [10](#), [13](#), [14](#)
- [9] 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. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. [2](#)
- [10] Blender Online Community. *Blender - a 3D modelling and rendering package*. Blender Foundation, Stichting Blender Foundation, Amsterdam, 2018. [6](#), [10](#), [11](#)
- [11] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2017. [5](#)
- [12] Abe Davis, Marc Levoy, and Fredo Durand. Unstructured light fields. In *Computer Graphics Forum*, volume 31. Wiley Online Library, 2012. [2](#)
- [13] Nianchen Deng, Zhenyi He, Jiannan Ye, Budmonde Duinkharjav, Praneeth Chakravarthula, Xubo Yang, and Qi Sun. Fov-nerf: Foveated neural radiance fields for virtual reality. *IEEE Transactions on Visualization and Computer Graphics*, 28(11), 2022. [1](#)
- [14] Terrance DeVries, Miguel Angel Bautista, Nitish Srivastava, Graham W Taylor, and Joshua M Susskind. Unconstrained scene generation with locally conditioned radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. [3](#)
- [16] Yilun Du, Yanan Zhang, Hong-Xing Yu, Joshua B Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV) (ICCV)*. IEEE Computer Society, 2021. [3](#)
- [17] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. [4](#)
- [18] Hugging Face. *Stable Diffusion Demo*. [10](#)
- [19] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [20] Stephan J Garbin, Marek Kowalski, Matthew Johnson, Jamie Shotton, and Julien Valentin. Fastnerf: High-fidelity neural rendering at 200fps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. [3](#)
- [21] 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, 1996. 2

- [22] Karol Gregor, George Papamakarios, Frederic Besse, Lars Buesing, and Theophane Weber. Temporal difference variational auto-encoder. *arXiv preprint arXiv:1806.03107*, 2018. 3
- [23] Yuan-Chen Guo, Di Kang, Linchao Bao, Yu He, and Song-Hai Zhang. Nerfren: Neural radiance fields with reflections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [24] HomeTourVision. Real estate video tour — 7 rutledge ave, highland mills, ny 10930 — orange county, ny. 11
- [25] Tao Hu, Shu Liu, Yilun Chen, Tiancheng Shen, and Jiaya Jia. Efficientnerf efficient neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [26] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2, 3, 4
- [27] James T Kajiya and Brian P Von Herzen. Ray tracing volume densities. *ACM SIGGRAPH computer graphics*, 18(3), 1984. 3
- [28] Takuhiro Kaneko. Ar-nerf: Unsupervised learning of depth and defocus effects from natural images with aperture rendering neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 1, 2
- [29] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5, 11
- [30] 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. 5
- [31] Marc Levoy and Pat Hanrahan. Light field rendering. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, 1996. 2
- [32] Chaojian Li, Sixu Li, Yang Zhao, Wenbo Zhu, and Yingyan Lin. Rt-nerf: Real-time on-device neural radiance fields towards immersive ar/vr rendering. In *Proceedings of the 41st IEEE/ACM International Conference on Computer-Aided Design*, 2022. 1
- [33] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [34] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020. 3
- [35] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [36] Ricardo Martin-Brualla, Noha Radwan, Mehdi SM Sajjadi, Jonathan T Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [37] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (CVPR)*, June 2021. 3
- [38] 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)*. Springer, 2020. 1, 2, 5, 6, 7, 11
- [39] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4), 2022. 5, 6, 7, 10, 13, 14
- [40] Michael Oechsle, Songyou Peng, and Andreas Geiger. Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2
- [41] Onur Özyeşil, Vladislav Voroninski, Ronen Basri, and Amit Singer. A survey of structure from motion\*. *Acta Numerica*, 26, 2017. 2
- [42] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo Martin-Brualla. Nerfies: Deformable neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2, 3
- [43] Keunhong Park, Utkarsh Sinha, Peter Hedman, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Ricardo Martin-Brualla, and Steven M Seitz. Hypernerf: A higher-dimensional representation for topolog-

- ically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. 2, 3
- [44] Jialun Peng, Dong Liu, Songcen Xu, and Houqiang Li. Generating diverse structure for image inpainting with hierarchical vq-vae. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [45] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [46] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in Neural Information Processing Systems (NeurIPS)*, 32, 2019. 3
- [47] Sotheby’s International Realty. Spanish colonial retreat in scottsdale, arizona. 11
- [48] Daniel Rebain, Mark Matthews, Kwang Moo Yi, Dmitry Lagun, and Andrea Tagliasacchi. Lolnerf: Learn from one look. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3
- [49] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 3
- [50] Konstantinos Rematas, Andrew Liu, Pratul P Srinivasan, Jonathan T Barron, Andrea Tagliasacchi, Thomas Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [51] Mehdi SM Sajjadi, Henning Meyer, Etienne Pot, Urs Bergmann, Klaus Greff, Noha Radwan, Suhani Vora, Mario Lučić, Daniel Duckworth, Alexey Dosovitskiy, et al. Scene representation transformer: Geometry-free novel view synthesis through set-latent scene representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [52] Johannes L Schonberger and Jan-Michael Frahm. Structure-from-motion revisited. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 1, 2
- [53] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020. 3
- [54] Harry Shum and Sing Bing Kang. Review of image-based rendering techniques. In *Visual Communications and Image Processing 2000*, volume 4067. SPIE, 2000. 2
- [55] Pratul P Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [56] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [57] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Blocknerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [58] Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in Neural Information Processing Systems (NeurIPS)*, 33, 2020. 3
- [59] Bayer Video Tours. 31 brian dr, rochester, ny presented by bayer video tours. 11
- [60] Edgar Tretschk, Ayush Tewari, Vladislav Golyanik, Michael Zollhöfer, Christoph Lassner, and Christian Theobalt. Non-rigid neural radiance fields: Reconstruction and novel view synthesis of a dynamic scene from monocular video. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 3
- [61] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [62] Aaron Van Den Oord, Oriol Vinyals, et al. Neural discrete representation learning. *Advances in Neural Information Processing Systems (NeurIPS)*, 30, 2017. 2, 3, 4
- [63] Dor Verbin, Peter Hedman, Ben Mildenhall, Todd Zickler, Jonathan T Barron, and Pratul P Srinivasan. Ref-nerf: Structured view-dependent appearance for neural radiance fields. In *2022 IEEE/CVF Conference*

- on *Computer Vision and Pattern Recognition (CVPR) (CVPR)*. IEEE, 2022. 3
- [64] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yanshun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plenotrees for dynamic radiance field rendering in real-time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [65] Xin Wang, Shinji Takaki, Junichi Yamagishi, Simon King, and Keiichi Tokuda. A vector quantized variational autoencoder (vq-vae) autoregressive neural  $f_0$  model for statistical parametric speech synthesis. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28, 2019. 3
- [66] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4), 2004. 7
- [67] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 2
- [68] Liwen Wu, Jae Yong Lee, Anand Bhattad, Yu-Xiong Wang, and David Forsyth. Diver: Real-time and accurate neural radiance fields with deterministic integration for volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2
- [69] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [70] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qin-hong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021. 5, 6, 7
- [71] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021. 3
- [72] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2, 3
- [73] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018. 7
- [74] Tinghui Zhou, Richard Tucker, John Flynn, Graham Fyffe, and Noah Snavely. Stereo magnification: Learning view synthesis using multiplane images. *arXiv preprint arXiv:1805.09817*, 2018. 10, 11, 12, 13, 14