

Mixed Neural Voxels for Fast Multi-view Video Synthesis

Feng Wang¹ Sinan Tan¹ Xinghang Li¹ Zeyue Tian² Yafei Song³ Huaping Liu¹*

¹Beijing National Research Center for Information Science and Technology(BNRist),
Department of Computer Science and Technology, Tsinghua University

²Hong Kong University of Science and Technology

³ XR Lab, DAMO Academy, Alibaba Group

wang-f20@mails.tsinghua.edu.cn, hpliu@tsinghua.edu.cn

Abstract

Synthesizing high-fidelity videos from real-world multi-view input is challenging due to the complexities of real-world environments and high-dynamic movements. Previous works based on neural radiance fields have demonstrated high-quality reconstructions of dynamic scenes. However, training such models on real-world scenes is time-consuming, usually taking days or weeks. In this paper, we present a novel method named MixVoxels to efficiently represent dynamic scenes, enabling fast training and rendering speed. The proposed MixVoxels represents the 4D dynamic scenes as a mixture of static and dynamic voxels and processes them with different networks. In this way, the computation of the required modalities for static voxels can be processed by a lightweight model, which essentially reduces the amount of computation as many daily dynamic scenes are dominated by static backgrounds. To distinguish the two kinds of voxels, we propose a novel variation field to estimate the temporal variance of each voxel. For the dynamic representations, we design an inner product time query method to efficiently query multiple time steps, which is essential to recover the high-dynamic movements. As a result, with 15 minutes of training for dynamic scenes with inputs of 300-frame videos, MixVoxels achieves better PSNR than previous methods. For rendering, MixVoxels can render a novel view video with 1K resolution at 37 fps. Codes and trained models are available at <https://github.com/fengres/mixvoxels>.

1. Introduction

Dynamic scene reconstruction from multi-view videos is a critical and challenging problem, with many potential applications such as interactively free-viewpoint control

*Corresponding author.

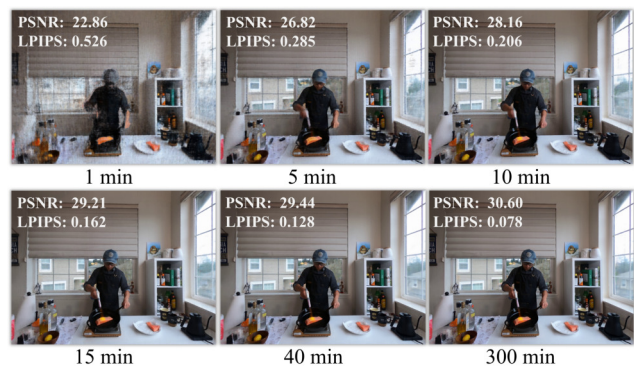


Figure 1. Our method enables rapid reconstruction of 4D dynamic scenes. We visualize the rendering results with different training schedules. With only 15 minutes of training, our approach achieves comparable PSNRs to other methods. Increasing the training time further enhances the ability to recover fine details.

for movies, cinematic effects like freeze-frame *bullet time*, novel view replays for sporting events, and various potential VR/AR applications. Recently, neural radiance fields [26] have demonstrated the possibility of rendering photo-realistic novel views for static scenes, with physically motivated 3D density and radiance modelling. Many methods [19, 20, 49, 13, 10, 31, 28, 29, 44] extend the neural radiance fields to dynamic scenes with additional time queries or an explicit deformation field. Many of these methods focus on the monocular input video setting on relatively simple dynamic scenes. To model more complex real-world dynamic scenes, a more practical solution is to use multi-view synchronized videos to provide dense spatial-temporal supervisions [55, 23, 3, 19].

Recently, Li et al. [19] propose a real-world dynamic scene dataset including many challenging situations such as objects of high specularities, topology changes, and volumetric effects. They address the problem by a hierarchical training scheme and the ray importance sampling strategies. Although significant improvements have been achieved, some challenges still exist: (1) The training and rendering take a

lot of time and computation resources. (2) Highly dynamic scenes with complex motions are still difficult to track.

In this paper, we focus on the multi-view 3D video synthesis problem and present a novel method named MixVoxels to address the above two challenges. The proposed MixVoxels is based on the explicit voxel-grid representation, which is recently popular due to its fast training and rendering speed on static scenes [50, 40, 8, 27]. We extend the voxel-grid representations to support dynamic scenes and propose an efficient inner product time querying method that can query a large number of time steps simultaneously, which is essential to recover the sharp details for highly-dynamic objects. Additionally, we represent dynamic scenes as a mixed static-dynamic voxel-grid representation. Specifically, the 3D spaces are split into static and dynamic voxels by our proposed variation field. The two components are processed by different models to reduce the redundant computations for the static space. Theoretically, once a dynamic scene consists of some static spaces, the training speed will benefit from the proposed mixed voxels. For a variety of events that occur in the physical world, the static components of environments are dominated in most cases, and the mixed voxels will speed up the training significantly in these scenarios. Besides, the separation of voxels makes the time-variant model focus on the dynamic regions, avoiding the time-aware voxels being biased by the static spaces to produce blurred motions. Our empirical validation confirms that the separation enables the model to learn sharp and distinct boundaries in high-dynamic regions. This also frees our method from the complex importance sampling strategies. With these designs, our method is capable of reconstructing a dynamic scene consisting of 300 frames within 15 minutes. To summarize, the main contributions of this work are:

- We propose a simple yet effective dynamic representation with inner product time querying method that can efficiently query multiple times simultaneously, improving the rendering quality for dynamic objects.
- We design an efficient variation field to separate static and dynamic spaces and present a mixed voxel-grid representation to accelerate training and rendering.
- We conduct qualitative and quantitative experiments to validate our method. As a result, the proposed MixVoxels achieves competitive or better rendering qualities with a $5000\times$ training speedup compared to implicit dynamic scene representations.

2. Related Works

Novel View Synthesis for Static Scenes. Synthesizing novel views for static scenes is a classical and well-studied problem. Different approaches represent the underlying geometric with different representations. Mesh-based meth-

ods [6, 9, 46, 48, 34, 43] represent the scenes with surfaces which is compact and easy to render, while optimizing a mesh to fit complex scenes is challenging. Volume-based methods such as voxel-grid [17, 35, 30, 23, 36] and multi-plane images (MPIs) [54, 11, 25, 39, 38, 45] are more suitable to model the complex and translucent scenes such as smooth and fluid. Particularly, Neural radiance fields [26] represent the scenes with an implicit volumetric neural representation, which employs a coordinate-based neural network to query the density and color for each point. The achieved photo-realistic rendering quality of NeRF led to an explosion of developments in the field. Advances have been made including improving the rendering qualities [42, 4], adapting to more general scenarios [52, 24, 41, 5], accelerating rendering or training speed [22, 51, 50, 40, 8], etc.

Novel View Synthesis for Dynamic Scenes. Synthesizing novel views for dynamic scenes is a more challenging and applicable problem. Recently, many extensions of NeRF for non-rigid dynamic scenes were proposed, which take a monocular video as input to learn the deformation and radiance fields. These methods can be categorized to modelling deformation implicitly [20, 49, 13, 10] (learn the non-decoupled deformation and appearance jointly) and explicitly [31, 28, 29, 44] (learn separated deformation and radiance fields and the deformation fields are usually in the form of relative motion with a canonical static space). Though improvements are achieved, reconstructing the complex general scenes is still difficult with only monocular videos. Most methods are constrained to fixed scenes like human-model or restricted motions. For real-world complex scenes, reconstructing from synchronized multi-view videos is more promising due to the dense supervision for every viewpoint and time instant. Earlier works [14, 55] explore the problem and show the possibility of rendering novel videos from a set of input views. Neural Volumes [23] proposes to use volumetric representations. They employ an encoder-decoder network to convert input images into a 3D volume, and decode the latent representations by the differentiable ray marching operation. [3] presents a data-driven approach for 4D space-time visualization of dynamic scenes by splitting static and dynamic components and using a U-Net structure in screen space to convert intermediate representation to image. Different with this method, our method split the static and dynamic components in the 3D voxel space instead of the pixel space. More recently, DyNerf [19] uses a temporal-aware neural radiance field to address the problem, and proposes some sampling strategies to train it efficiently. Compared with previous methods, they propose a more complicated real-world dataset and validate their method. For accelerating the reconstruction of dynamic scenes, FourierPlenotree [47] proposes to model the dynamics in frequency domain, and generate a Plenotree through multi-

view blending to accelerate rendering. They focus on the foreground moving objects extracted via chroma key segmentation, which requires the background should be a pure color (or rely on segmentation algorithms). Recently, the acceleration of training and rendering for dynamic scenes has attracted much attention. Cocurrent works include StreamRF [18] which proposes to accelerate the training of dynamic scenes by modeling the differences of adjacent frames, NeRFPlayer [37] which decomposes the dynamic scenes into static, new and deforming components, Hyperreel [2] which proposes an efficient sampling network and models keyframes. K-Planes [12] and HexPlanes [7] decompose the 4D dynamic scenes into different 2D representations.

Acceleration of Neural Radiance Fields. While Neural radiance fields can render novel views with high fidelity, training and rendering require querying a deep MLP millions of times which is computationally intensive. Many recent methods propose to accelerate the training and rendering speed of NeRF. For rendering, Neural Sparse Voxel Fields [22] proposes a voxel-grid representation to skip over many empty regions. PlenOctree [51] accelerates the rendering process by pre-tabulating the NeRF into a PlenOctree and using the spherical harmonic representation of radiance. Derf [32] and Kilonerf [33] propose to accelerate the rendering speed by dividing the scenes into multiple areas, and employ multiple small network in each area. AutoInt [21] proposes to restructure the MLP network to accelerate the computations of ray integrals, which helps accelerate the rendering speed. For accelerating the training of NeRF, some methods use explicit voxel-grid representations [50, 40] to accelerate the training process and convergence speed. Instant-NGP [27] proposes a multi-resolution hash table structure to accelerate the training. The model sizes of most fast training methods are relatively large due to a large number of voxels. TensoRF [8] proposes to reduce the model size by factorizing the 4D scene tensor into multiple compact low-rank tensor components.

3. Method

In this section, we introduce the proposed MixVoxels, which represents the 4D dynamic scenes as mixtures of static and dynamic voxels. Fig. 2 illustrates the overview of our method. In the following subsections, we will first introduce the voxel-grid representations for static scenes and our extension to dynamic scenes. Then we introduce the variation field for identifying the dynamic voxels. At last, we introduce the training of MixVoxels.

3.1. Static Voxel-grid Representation

Neural radiance fields [26] have demonstrated photo-realistic novel viewpoint synthesis, while the training of NeRF requires extensive computation due to millions of

neural network queries. For accelerating NeRF, many recent works [22, 50, 40, 8] have explored the explicit volumetric representation, which avoids the huge amount of computation of querying neural network. Specifically, a 3D scene is split into $N_x \times N_y \times N_z$ voxels. The densities and color features are stored in these voxels and denoted as $S^\sigma \in \mathbb{R}^{N_x \times N_y \times N_z}$ and $S^c \in \mathbb{R}^{N_x \times N_y \times N_z \times C}$. $S_{i,j,k}^\sigma$ and $S_{i,j,k}^c$ represent the learnable density and color feature of the voxel corner at a discrete position (i, j, k) . For a continuous position (x, y, z) , the representation $S_{x,y,z}$ can be calculated by interpolating the nearest 8 discrete positions. A small MLP network C_θ is used to parse the color features into RGB values, taking S^c and view direction \mathbf{d} as input. Formally, the density σ and color c is formulated as

$$\sigma(x, y, z) = S_{x,y,z}^\sigma, \quad c(x, y, z, \mathbf{d}) = C_\theta(S_{x,y,z}^c, \mathbf{d}). \quad (1)$$

3.2. Dynamic Voxel-grid Representation

For dynamic scenes, a direct extension is to add the time dimension to the static voxel-grid representation \mathcal{S} explicitly. However, this direct extension is almost memory-prohibitive due to the large and linearly increasing memory footprint. For a 300-frame video, the learned models will occupy 30 GB of memory and be difficult to train with GPUs due to the limitation of GPU memory. To address this problem, we propose a spatially explicit and temporally implicit representation to reduce the memory footprint. Specifically, we represent the dynamic scene as a 4D learnable voxel-grid $\mathcal{G}^\sigma \in \mathbb{R}^{N_x \times N_y \times N_z \times C_1}$ and $\mathcal{G}^c \in \mathbb{R}^{N_x \times N_y \times N_z \times C_2}$. Different from the static scene representations, the densities and colors for all time steps are implicitly encoded as compact features stored in each voxel corner. The compact features will be processed by a time-aware projection to acquire density and color for each time step. Concretely, for the compact density feature $\mathcal{G}_{x,y,z}^\sigma$ and color feature $\mathcal{G}_{x,y,z}^c$ in any position (x, y, z) , we employ two MLPs $\mathcal{T}_{\theta_1}^\sigma$ and $\mathcal{T}_{\theta_2}^c$ to increase the feature dimensions for better parsing time-variant density and color. The MLPs here can be viewed as decompressors that decompress the compact low-dimensional voxel-grid features into more tractable ones. Compared with directly storing high-dimensional features in each voxel, the temporally implicit representation reduces the memory footprint significantly since the shared MLPs only increase memory slightly.

Inner product time query. For a discrete time step t , we use a learnable time-variant latent representation ω_t to represent the time query. Instead of concatenating the time query with the intermediate features, we propose to calculate the inner product between the learned time query and the decompressed features as the required output σ and c . Formally, the density and color of a space-time query (x, y, z, t) are formulated as

$$\sigma(x, y, z, t) = \omega_t^\sigma \cdot \mathcal{T}_{\theta_1}^\sigma(\mathcal{G}_{x,y,z}^\sigma), \quad (2)$$

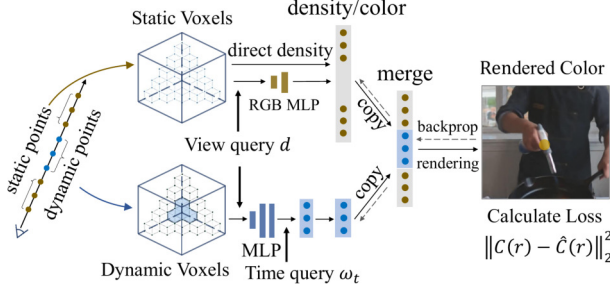


Figure 2. Overview of our method. Given a ray, we first sample points, and split them into static and dynamic ones using the variation field. After that, we feed these points to the corresponding branches and query the required properties. Then we merge the static output and dynamic output for rendering the ray color. An L2 loss is employed to calculate loss and back-propagate.

$$c(x, y, z, \mathbf{d}, t) = \omega_t^c \cdot \mathcal{T}_{\theta_2}^c(\mathcal{G}_{x,y,z}^c, \mathbf{d}). \quad (3)$$

In practice, simultaneously querying multiple time steps helps reconstruct the detail of high-dynamic motions and reduces the training iterations to traverse through all time steps. The inner product based query will facilitate the training speed when simultaneously querying many time steps in a training iteration. Specifically, we denote the FLOPs of the MLP \mathcal{T}_{θ_1} and the inner product operation as FLOP_{mlp} and FLOP_{inn} , respectively. For a T -frame video, the FLOPs of the concatenation query [19] is larger than $T \cdot \text{FLOP}_{mlp}$ (due to the extra temporal embedding dimension), while the FLOPs of our inner product query is only $\text{FLOP}_{mlp} + T \cdot \text{FLOP}_{inn}$ ($\text{FLOP}_{mlp} \gg \text{FLOP}_{inn}$).

3.3. Variation Field

In this subsection, we introduce the variation field to identify which voxels in the 3D space are dynamic, *i.e.*, the densities or colors are not constant over different time steps. By separating the static and dynamic voxels, the redundant computations caused by using a relatively heavy time-varying model to process the static components will be avoided, which accelerates the training and rendering.

A simpler solution for accelerate training is to separate the static and dynamic regions in pixel-level, *i.e.*, using the temporal variance of pixels to produce static and dynamic ones. However this scheme is actually not feasible because we can only separate dynamic and static regions in training views using the ground truth. For rendering novel views, we can not get the pixel variance for rendering since we have no ground truth in novel views. Thus a feasible solution is to learn the voxel-level temporal variance which is shared for all possible views. In addition, separating in the voxel-level is more efficient compared to pixel-level, even if we have an oracle to make the pixel-level separation feasible in novel views. This is because not all voxels projected to a dynamic pixel are dynamic, there will be only a small fraction of voxels around the object surfaces are actually

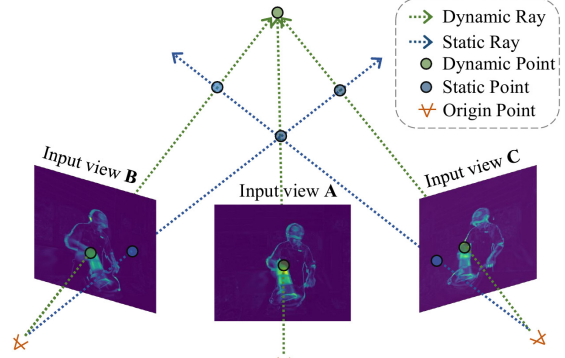


Figure 3. Implicit interaction of multiple rays to decide which point is dynamic. For a static ray (blue line), all points are set to be low dynamic. For a dynamic ray (green line), at least one point is dynamic. The intersection of multiple dynamic ray is more likely to be a dynamic point, which is also physically intuitive.

dynamic. Therefore, the voxel level separation will produce much fewer dynamic queries.

To perform the voxel-level separation, we utilize the pixel-level temporal variances from training videos as the supervision to estimate the voxel-level variances. The pixel-level (or ray-level) temporal variances of different videos are shown in Fig. 3. Formally, given a ray $\mathbf{r}(s) = \mathbf{o} + s \cdot \mathbf{d}$ with origin \mathbf{o} and direction \mathbf{d} , the corresponding pixel color at time step t is defined as $C(\mathbf{r}, t)$. Then the pixel-level temporal variance $D^2(\mathbf{r})$ is formalized as

$$D^2(\mathbf{r}) = \frac{1}{T} \sum_{t=1}^T (C(\mathbf{r}, t) - \bar{C}(\mathbf{r}))^2, \quad \bar{C}(\mathbf{r}) = \frac{1}{T} \sum_{t=1}^T C(\mathbf{r}, t),$$

where $\bar{C}(\mathbf{r})$ is the mean color of pixel corresponding to the ray \mathbf{r} . For identifying the dynamic pixels, the standard deviation $D(\mathbf{r})$ is binarized to $M(\mathbf{r})$ with a threshold γ to provide pixel-level dynamic supervision, *i.e.* $M(\mathbf{r}) = 1$ if $D(\mathbf{r}) \geq \gamma$, else $M(\mathbf{r}) = 0$. In this way, we judge that a ray \mathbf{r} is dynamic if $M(\mathbf{r}) = 1$. Next, we use the $M(\mathbf{r})$ as supervision to estimate the voxel-level variations \mathcal{V} .

The relations between the pixel-level variance and voxel-level variance lie in the following aspects: (1) If a pixel is static, then all voxels passed through by the ray corresponding to the pixel should be static in most cases (We will discuss some special cases which violate this rule later in this subsection.). (2) If a pixel is dynamic, then at least one of the voxels passed through by the corresponding ray is dynamic. Fig. 3 shows the two situations. With the above two relations, we design the variation field, which is denoted as $\mathcal{V} \in \mathbb{R}^{N_x \times N_y \times N_z}$ to represent the voxel-level temporal variance. Specifically, we uniformly sample N_s points from the near plane to the far plane in \mathbf{r} , and build the following equation to satisfy the two relations mentioned above:

$$\hat{M}(\mathbf{r}) = \mathbf{s}(\max(\{\mathcal{V}_{\mathbf{r}(s_i)} | i \in \{1, \dots, N_s\}\})), \quad (4)$$

where $\hat{M}(\mathbf{r})$ is an estimation of $M(\mathbf{r})$, and \mathbf{s} is the sigmoid function. Then we train the variation field by minimizing

the following binary cross-entropy loss:

$$\mathcal{L}_v = \mathbb{E}_r \left[-M(\mathbf{r})\log(\hat{M}(\mathbf{r})) - (1 - M(\mathbf{r}))\log(1 - \hat{M}(\mathbf{r})) \right]. \quad (5)$$

By optimizing the above loss function to all rays, we can get the learned variation field \mathcal{V} . The training of the variation field is very efficient, usually taking *less than 30 seconds*.

The maximization operation well formulates the relations between a pixel and its corresponding voxels. If a pixel is static, then the equations of Eq. (4) and Eq. (5) will force all voxels passed through by the corresponding ray to be static ($\mathcal{V}_{x,y,z} = 0$). If a pixel is dynamic, Eq. (4) requires at least one of the corresponding voxels (*i.e.*, the max value of the voxel variances) to be dynamic ($\mathcal{V}_{x,y,z} = +\infty$). Although we provide no information about which specific voxels in a dynamic ray are dynamic, the implicit interaction of multiple different rays will force the solution to be physically reasonable. To explain this, we focus on the observable voxels which at least passed through by one ray. If a point (x, y, z) is passed through by at least one static ray, then $\mathcal{V}_{x,y,z}$ will tend to be optimized to be close to zero. If a point (x, y, z) is only passed through by dynamic rays, and not occluded by other dynamic voxels, then $\mathcal{V}_{x,y,z}$ will be optimized to $+\infty$. This is because without occlusion, the above point of (x, y, z) along the dynamic rays will be passed by other static rays (the front space along these rays are observable from other views). We illustrate this situation in Fig. 3.

Inference. After the training process, the temporal variation at a specific 3D position (x, y, z) is $\mathcal{V}_{x,y,z}$, which is easily acquired by interpolating the discrete variation field. We then identify a voxel in the scene as dynamic if $\mathcal{V}_{x,y,z}$ is larger than a hyper-parameter β , and as static if it is smaller than β . Formally, we will get a dynamic mask $\hat{M} \in \{0, 1\}^{N_x \times N_y \times N_z}$, which will be used to split sampling points in a ray into static points and dynamic points. We evaluate the effectiveness of this inference method in the test views and find that the recall and precision are reasonable for splitting the dynamic and static parts (recall: 0.97, precision: 0.94 when $\beta = 0.9$). Although the recall seems sufficient to retrieve most dynamic parts, we empirically find some false negatives in the rendering images affect the rendering quality. To address this problem, we use a ray-wise max-pooling operation to identify the points near to a dynamic point as dynamic. The kernel size of max-pooling is set to $k_m = 21$, and the stride is set to 1. In this way, the recall is very closed to 1. Although many hyper-parameters are incorporated, we have empirically found that the thresholds γ and β are not sensitive over a wide range of reasonable values.

Discussion. There may be some situations in which the rules (1) are broken due to occlusion. Specifically, when a dynamic voxel is occluded by some static voxels, the occluded parts should not be classified as static, which makes the 2D supervision noisy. However in practice, we found



Figure 4. Impact of occlusions to the variation field.

the learning-based separation has a certain degree of tolerance for this situation. We conducted experiments to verify this and present the results in Fig. 4. The occluded region is visualized in the leftmost view (the area behind the static roadblock). Although the dynamic region marked in yellow is occluded in the left view, it is actually classified as dynamic region. This can be inferred from another view where the occluded region has changed, as illustrated in the middle and right images in Fig. 4). The variation field is learning-based and will learn a solution that satisfies most constraints. If it forces some dynamic voxels occluded by static voxels to be 0, then the loss function from other visible views will be high. As a result, the learning-based process tends to assign a “middle solution” to voxels with inconsistent supervisions from different views. We also attempted to use the transmittance as a weight (in a way of volumetric rendering) to learn the variation field which can explicitly handle the occlusion problem. However, we found a large efficiency drop with similar performance. As a result, we use the proposed variation field and find this formulation works well for most scenes, including some challenging scenes with large areas of motions.

3.4. Training of Mixed Neural Voxels

With the help of the variation field, we can split a scene into dynamic voxels and static voxels. To reduce redundant computations, we use the lightweight static model described in Sec. 3.1 to compute the densities and colors for static voxels and the dynamic model described in Sec. 3.2 to compute the densities and colors for dynamic voxels. The overall architecture is illustrated in Fig. 2.

Specifically for a given ray $\mathbf{r}(s) = \mathbf{o} + s\mathbf{d}$ with origin \mathbf{o} and view direction \mathbf{d} , we apply stratified sampling from the near to the far planes and get N_s points. Then the N_s points are separated into static and dynamic ones by inferring these points with the proposed variation field. For the static points, we pass them into the static branch to retrieve the colors and densities. For the dynamic points, we pass them into the dynamic branch together with a deferred time query ω_t to retrieve the corresponding properties. After that, we merge the static points and dynamic points according to their order. Then we apply volumetric rendering to the merged points to obtain the rendered color, which is formulated as

$$C(\mathbf{r}, t) = \sum_{i=1}^{N_s} T_{i,t} \cdot (1 - \exp(-\sigma_{i,t}\delta_i)) \cdot c_{i,t}, \quad (6)$$

where $T_{i,t}$ is the accumulated opacity (or transmittance):

$T_{i,t} = \exp(-\sum_{j=1}^{i-1} \sigma_{j,t} \delta_j)$, and δ_i is the distance between adjacent samples. Given the ground truth color $C_g(\mathbf{r}, t)$, an l_2 -loss is employed to train the model:

$$\mathcal{L} = \mathbb{E}_{(\mathbf{r}, t)} [\|C_g(\mathbf{r}, t) - C(\mathbf{r}, t)\|_2^2]. \quad (7)$$

For both static and dynamic branches, we omit the computation of color for points whose densities are close to zero, which is a widely adopted pruning strategy [22, 50, 8].

Efficiency analysis. We define the proportion of dynamic points in a scene as λ (≈ 0.05 for most scenes). Besides, the FLOPs of static and dynamic branches are denoted as FLOP_{sta} and FLOP_{dyn} , respectively. Then the total FLOPs of MixVoxels are $\text{FLOP}_{\text{sta}} + \lambda \cdot \text{FLOP}_{\text{dyn}}$. Empirically, $\text{FLOP}_{\text{dyn}} / \text{FLOP}_{\text{sta}}$ ranges from 50 to 100 with different reasonable dimension settings. Then the acceleration ratio of splitting static and dynamic models is $\text{FLOP}_{\text{dyn}} / (\text{FLOP}_{\text{sta}} + \lambda \text{FLOP}_{\text{dyn}}) \approx 10$. In practice, the actual speedup with a 3090 GPU is about 5, the inconsistency between analysis and experiment may come from the GPU features, which are friendly to a more consistent network.

3.5. Implementation Detail

The voxel-grid representations require large GPU memory to store the cubically growing voxel numbers. To implement the voxel-grid representations more memory efficient, we use the tensor factorization technique proposed in TensorRF [8] to reduce the memory footprint. In this way, a 3D tensor is factorized into the outer product of a vector and a 2D matrix. We factorize all the voxel-grid tensors, including static voxels, dynamic voxels and the variation field. With the help of tensor factorization, the learned model costs about 500MB for a 300-frame multi-view video scene. For the voxel resolutions, we follow [8] to start from an initial low resolution of 256^3 , and upsample the resolution at steps 1500, 2000, 2500, and 2750 with a linear increase in the log space. The final resolution is set to 640^3 . Once the resolution is changed, we re-train the variation field, which only takes about 15-30s. The voxel-grid feature dimension is set to 27, and the hidden state of MLP is set to 512. For training, we use Adam [15] optimizer with a learning rate of 0.02 for voxels and $3e-3$ for MLPs. The total variation loss [50] is incorporated as a regularization to encourage the space smoothness.

4. Experiments

4.1. Experiment Setting

Dataset. We validate our method on two datasets: (1) The Plenoptic Video Dataset [19], which consists of 6 publicly accessible scenes: coffee-martini, flame-salmon, cook-spinach, cut-roasted-beef, flame-steak and sear-steak. We conduct experiments on all six scenes. Each scene contains 19 videos with different camera views. The dataset

Table 1. Results on our collected dataset, including two scenes.

| Scene | Model | PSNR \uparrow | DSSIM \downarrow | LPIPS \downarrow |
|----------------------|-------------|-----------------|--------------------|--------------------|
| Moving-Cars | MixVoxels-S | 18.72 | 0.251 | 0.689 |
| | MixVoxels-M | 18.97 | 0.228 | 0.552 |
| | MixVoxels-L | 18.89 | 0.222 | 0.540 |
| | MixVoxels-X | 19.11 | 0.210 | 0.516 |
| Solving-Rubik | MixVoxels-S | 25.39 | 0.065 | 0.339 |
| | MixVoxels-M | 26.05 | 0.059 | 0.275 |
| | MixVoxels-L | 26.28 | 0.055 | 0.241 |
| | MixVoxels-X | 26.80 | 0.047 | 0.209 |

contains many challenging scenes including objects with topology changes, objects with volumetric effects, various lighting conditions, etc. (2) Our proposed dataset including two more complex dynamic scenes: moving-cars and solving-rubik. The moving-cars scene features several vehicles passing across the screen, with significant motion and displacement. Meanwhile, in the solving-rubik scene, a man solves a Rubik’s cube at a rapid pace, averaging 4 rotations per second, providing an opportunity to evaluate the model’s ability to capture swift movements. The collection procedures used are similar to those of DyNeRF. More details are presented in the appendix.

For training and evaluation, we follow the experiment setting in [19] that employs 18 views for training and 1 view for evaluation. To quantitatively evaluate the rendering quality on novel views, we measure PSNR, DSSIM and LPIPS[53] on the test views. We also provide more metrics in the appendix including FLIP [1] and JOD [16], which we find the comparisons are similar with PSNR and LPIPS. We follow the setting of [19] to evaluate our model frame by frame. For videos consisting of equal or more than 300 frames, we evaluate our model every 10 frames [19] to calculate the frame-by-frame metrics except for the JOD metrics, which requires a stack of continuous video.

Training Schedules. For evaluating the effect of training time, we train MixVoxels with different configurations shown in Tab. 2. The configurations vary in terms of training iterations and the number of sample points per ray. By default, the step size for sampling points is set to four times of the voxel width. The $8\times$ means that there will be eight times as many sampling points compared to the default.

Table 2. Different training configurations of MixVoxels.

| Model | Iterations | Sampling points | Training Time |
|-------------|------------|-----------------|---------------|
| MixVoxels-S | 5000 | $1\times$ | 15 min |
| MixVoxels-M | 12500 | $1\times$ | 40 min |
| MixVoxels-L | 25000 | $1\times$ | 80 min |
| MixVoxels-X | 50000 | $8\times$ | 300 min |

4.2. Results

Quantitative results and comparisons. For quantitative results, we present the metrics and compare with other methods in Tab. 3. Compared with the previous state-of-the-art method DyNeRF, we reduce the training time from 1.3K GPU hours to 15 minutes, making the training of complex



Figure 5. Visual comparisons with state-of-the-art methods. K-Planes [12] and HexPlane [7] are concurrent works. We have selected four representative patches to better inspect the details. Our method performs well on reconstructing details and capturing movements.

Table 3. Quantitative results comparisons. All metrics are measured on 300-frame scenes. We also report the training time, rendering speed (FPS) and model size. * Note DyNeRF is trained on 8 GPUs, while others are trained on one GPU.

| Method | Train | Render | Size | PSNR \uparrow | DSSIM \downarrow | LPIPS \downarrow |
|-----------------|---------|--------|--------|-----------------|--------------------|--------------------|
| DyNeRF[19] | 7 days* | - | 28 MB | 29.58 | 0.0197 | 0.083 |
| Concurrent work | | | | | | |
| StreamRF[18] | 75 min | 8.3 | 5310MB | 28.26 | - | - |
| NeRFPlayer[37] | 360 min | 0.05 | - | 30.69 | 0.034 | 0.111 |
| Hyperreel[2] | 540 min | 2.0 | 360 MB | 31.10 | 0.036 | 0.096 |
| K-Planes[12] | 108 min | - | - | 31.63 | 0.018 | - |
| HexPlanes[7] | 720 min | - | 200MB | 31.71 | 0.014 | 0.075 |
| MixVoxels-S | 15 min | 37.7 | 500 MB | 31.03 | 0.022 | 0.129 |
| MixVoxels-M | 40 min | 37.7 | 500 MB | 31.22 | 0.019 | 0.102 |
| MixVoxels-L | 80 min | 37.7 | 500 MB | 31.34 | 0.017 | 0.096 |
| MixVoxels-X | 300 min | 4.6 | 500 MB | 31.73 | 0.015 | 0.064 |

dynamic scenes more practical. For rendering, the MixVoxels has a 37.7 fps rendering speed for 1K resolution. Compared with concurrent works, MixVoxels requires less training time and achieves faster rendering speed, while achieving competitive PSNR and LPIPS. For example, with only 15 minutes of training, MixVoxels achieve 31.03 PSNR which is comparable to other methods trained for hours. With sufficient training, all metrics are further improved. For the quantitative results on our collected more complex scenes, we present them on Tab. 1.

Qualitative results and comparisons. Fig. 6 demonstrates the novel view rendering results on different dynamic scenes. The first four rows are novel view videos from Plenoptic Video Dataset [19]. The last two rows present the novel view videos from our collected two more complex dynamic scenes. The results show that our method can achieve near photo-realistic rendering quality. We provide the video results at the supplementary material. For qualitative comparisons, we show them in Fig. 5. MixVoxels can better reconstruct the moving object (the firing gun) and textual details like the hat and the salmon stripes.

We further investigate the relations between rendering efficiency and rendering quality. As shown in the lower part of Tab. 3, it was observed that an increase in training time leads to improvements in both PSNR and LPIPS. Longer

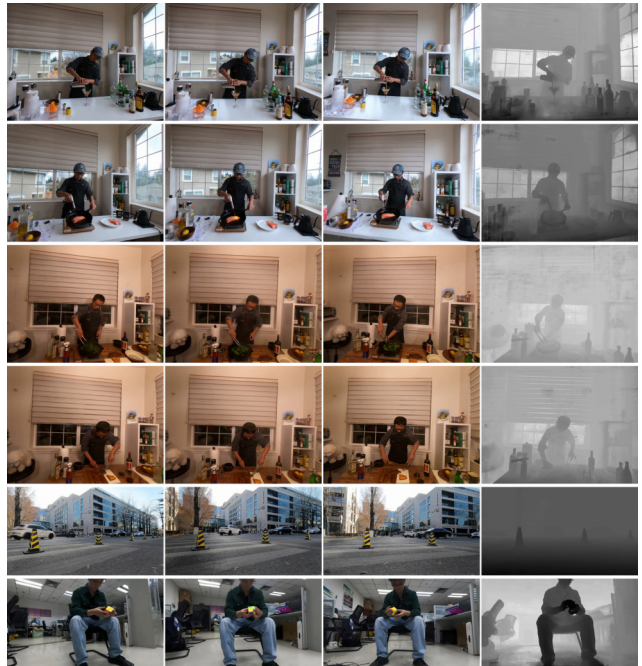


Figure 6. Novel view synthesis of MixVoxels. We select some frames at different views. The last column demonstrates the normalized depth. We provide videos at the supplemental material.

training facilitates the reconstruction of sharp boundaries and fine details. Visual comparisons presented in Fig. 7 reveals that 15 minutes of training produces satisfactory recovery of most scene components but resulted in blurry motion details. With longer training, the moving objects become clearer with a distinct boundary.

4.3. Ablation Study

In this subsection, we empirically justify the design of MixVoxels by ablating or modifying several key features. We also provide analysis that intuitively explains the ablations. We conduct all experiments in this subsection on the coffee-martini scene, which we find is typical for demonstrating the fast-moving complex objects.

Ablation on splitting voxels. To study the effect of splitting static and dynamic voxels, we compare MixVox-



Figure 7. Qualitative demonstration of different training schedules. Longer training helps better reconstruct the high-dynamic parts.



Figure 8. Qualitative comparison of MixVoxels and the full-dynamic model. Training with full-dynamic models with the same iterations can not well reconstruct the motion details.

Table 4. Ablation on the mixed voxels. Training with the full-dynamic voxel model hurts both efficiency and efficacy.

| Method | Time | PSNR \uparrow | DSSIM \downarrow | LPIPS \downarrow | FLIP \downarrow | JOD \uparrow |
|--------------|-------------|-----------------|--------------------|--------------------|-------------------|----------------|
| Full-Dynamic | 2.5h | 28.36 | 0.036 | 0.2236 | 0.1196 | 7.44 |
| MixVoxels | 0.6h | 29.47 | 0.026 | 0.1167 | 0.1223 | 7.99 |

Table 5. Ablation on three different methods for time query. We only substitute the time query with different method, and train them on the proposed MixVoxels framework.

| Method | Time | PSNR \uparrow | DSSIM \downarrow | LPIPS \downarrow | FLIP \downarrow | JOD \uparrow |
|---------------|------|-----------------|--------------------|--------------------|-------------------|----------------|
| Concat | 58m | 28.95 | 0.037 | 0.2146 | 0.1294 | 7.44 |
| Fourier | 43m | 28.67 | 0.029 | 0.1824 | 0.1286 | 7.60 |
| Inner product | 40m | 29.47 | 0.026 | 0.1167 | 0.1223 | 7.99 |

els with a full-dynamic voxel-grid representation, where all points are processed by the dynamic model. Tab. 4 shows the comparisons. With the same training iterations, the full-dynamic model is more time-consuming, which is intuitive because it processes all voxels with the dynamic models. Fig. 8 shows the qualitative comparison. The full-dynamic model recovered blurred motions. We speculate the reason is because the large area of static regions affects the capturing of dynamic information. The network will be biased by most static voxels with no motions and tend to learn low-frequency information.

Ablation on time query. We compare our inner product time query method with other variants: (1) Concatenation which concatenates the temporal embedding with the voxel features to be processed by an MLP. (2) Fourier head proposed by [47] which reconstructs the dynamics in



Figure 9. Ablation on the number of time query denoted as Q.

frequency-domain. Tab. 5 shows the performance comparison. The concatenation query method is both space- and time-consuming. Querying one time step requires forwarding the fused features through the whole MLP. Limited by the GPU memory, we can only query 50 time steps per iteration with the concatenation way, which harms the performance on high-dynamic regions. The Fourier head processes the features to predict the magnitudes of different frequency components and the performance is competitive, while it requires an additional inverse discrete Fourier transform to recover the information in the temporal domain. Overall, the inner product query is the simplest and most efficient way for querying.

Number of time queries per-iteration. We empirically find that simultaneously querying multiple time steps in an iteration helps reconstruct the details of moving parts. Fig. 9 demonstrates the effect of different numbers of time queries denoted as Q. With more time queries, the boundaries of the moving hand and the flowing coffee become clearer. Querying more time steps can provide dense supervision and make the model acquire global temporal information in every iteration, which accelerates the convergence speed. The effective inner product time queries allow adding more time queries with negligible increase in computation.

4.4. Limitations

Our method can synthesize novel view videos with a relative high quality. However, for some scenes with complex lighting conditions, some inconsistent property predictions may appear at the boundary between dynamic and static voxels, which is shown in Fig. 10. We suspect that the phenomenon is caused by the under-sampling of dynamic regions on scenes with some bad conditions. We will investigate ways to address the problem in future works.



Figure 10. Some inconsistent density and color predictions in the boundaries between dynamic and static regions.

5. Conclusion

This paper demonstrates a new method named MixVoxels to efficiently reconstruct the 4D dynamic scenes and synthesize novel view videos. The core of our method is to split the 3D space into static and dynamic components

with the proposed variation field, and process them with different branches. The separation speeds up the training and makes the dynamic branch focus on the dynamic parts to improve the performance. We also design an efficient dynamic voxel-grid representation with an inner product time query. The proposed method achieves competitive results with only 15 minutes of training, making the training and rendering of complex dynamic scenes more practical. We believe the fast training speed will enable potentially useful applications that are bottlenecked by training efficiency.

6. Acknowledgement

This work was supported in part by the National Natural Science Fund for Distinguished Young Scholars under Grant 62025304.

References

- [1] Pontus Andersson, Jim Nilsson, Tomas Akenine-Möller, Magnus Oskarsson, Kalle Åström, and Mark D Fairchild. Flip: A difference evaluator for alternating images. *Proc. ACM Comput. Graph. Interact. Tech.*, 3(2):15–1, 2020. 6
- [2] Benjamin Attal, Jia-Bin Huang, Christian Richardt, Michael Zollhoefer, Johannes Kopf, Matthew O’Toole, and Changil Kim. Hyperreel: High-fidelity 6-dof video with ray-conditioned sampling. *arXiv preprint arXiv:2301.02238*, 2023. 3, 7
- [3] Aayush Bansal, Minh Vo, Yaser Sheikh, Deva Ramanan, and Srinivasa Narasimhan. 4d visualization of dynamic events from unconstrained multi-view videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5366–5375, 2020. 1, 2
- [4] 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*, pages 5855–5864, 2021. 2
- [5] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 2
- [6] Chris Buehler, Michael Bosse, Leonard McMillan, Steven Gortler, and Michael Cohen. Unstructured lumigraph rendering. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 425–432, 2001. 2
- [7] Ang Cao and Justin Johnson. Hexplane: a fast representation for dynamic scenes. *arXiv preprint arXiv:2301.09632*, 2023. 3, 7
- [8] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *arXiv preprint arXiv:2203.09517*, 2022. 2, 3, 6
- [9] Paul E Debevec, Camillo J Taylor, and Jitendra Malik. Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 11–20, 1996. 2
- [10] 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)*, pages 14304–14314. IEEE Computer Society, 2021. 1, 2
- [11] John Flynn, Michael Broxton, Paul Debevec, Matthew Duvall, Graham Fyffe, Ryan Overbeck, Noah Snavely, and Richard Tucker. Deepview: View synthesis with learned gradient descent. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2367–2376, 2019. 2
- [12] Sara Fridovich-Keil, Giacomo Meanti, Frederik Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. *arXiv preprint arXiv:2301.10241*, 2023. 3, 7
- [13] 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*, pages 5712–5721, 2021. 1, 2
- [14] Takeo Kanade, Peter Rander, and PJ Narayanan. Virtualized reality: Constructing virtual worlds from real scenes. *IEEE multimedia*, 4(1):34–47, 1997. 2
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [16] Vamsi Kiran Adhikarla, Marek Vinkler, Denis Sumin, Rafal K Mantiuk, Karol Myszkowski, Hans-Peter Seidel, and Piotr Didyk. Towards a quality metric for dense light fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 58–67, 2017. 6
- [17] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *International journal of computer vision*, 38(3):199–218, 2000. 2
- [18] Lingzhi Li, Zhen Shen, Zhongshu Wang, Li Shen, and Ping Tan. Streaming radiance fields for 3d video synthesis. *arXiv preprint arXiv:2210.14831*, 2022. 3, 7
- [19] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, et al. Neural 3d video synthesis from multi-view video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5521–5531, 2022. 1, 2, 4, 6, 7
- [20] 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*, pages 6498–6508, 2021. 1, 2
- [21] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14556–14565, 2021. 3
- [22] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances*

- in *Neural Information Processing Systems*, 33:15651–15663, 2020. 2, 3, 6
- [23] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *arXiv preprint arXiv:1906.07751*, 2019. 1, 2
- [24] 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*, pages 7210–7219, 2021. 2
- [25] Ben Mildenhall, Pratul P Srinivasan, Rodrigo Ortiz-Cayon, Nima Khademi Kalantari, Ravi Ramamoorthi, Ren Ng, and Abhishek Kar. Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (TOG)*, 38(4):1–14, 2019. 2
- [26] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM*, 65(1):99–106, 2021. 1, 2, 3
- [27] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 2, 3
- [28] 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*, pages 5865–5874, 2021. 1, 2
- [29] 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 topologically varying neural radiance fields. *arXiv preprint arXiv:2106.13228*, 2021. 1, 2
- [30] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM Transactions on Graphics (TOG)*, 36(6):1–11, 2017. 2
- [31] 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*, pages 10318–10327, 2021. 1, 2
- [32] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14153–14161, 2021. 3
- [33] 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*, pages 14335–14345, 2021. 3
- [34] Gernot Riegler and Vladlen Koltun. Free view synthesis. In *European Conference on Computer Vision*, pages 623–640. Springer, 2020. 2
- [35] Steven M Seitz and Charles R Dyer. Photorealistic scene reconstruction by voxel coloring. *International Journal of Computer Vision*, 35(2):151–173, 1999. 2
- [36] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhofer. Deepvoxels: Learning persistent 3d feature embeddings. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2437–2446, 2019. 2
- [37] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerfplayer: A streamable dynamic scene representation with decomposed neural radiance fields. *arXiv preprint arXiv:2210.15947*, 2022. 3, 7
- [38] Pratul P Srinivasan, Ben Mildenhall, Matthew Tancik, Jonathan T Barron, Richard Tucker, and Noah Snavely. Lighthouse: Predicting lighting volumes for spatially-coherent illumination. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8080–8089, 2020. 2
- [39] Pratul P Srinivasan, Richard Tucker, Jonathan T Barron, Ravi Ramamoorthi, Ren Ng, and Noah Snavely. Pushing the boundaries of view extrapolation with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 175–184, 2019. 2
- [40] 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*, pages 5459–5469, 2022. 2, 3
- [41] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 2
- [42] 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*, 33:7537–7547, 2020. 2
- [43] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2
- [44] 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*, pages 12959–12970, 2021. 1, 2
- [45] Richard Tucker and Noah Snavely. Single-view view synthesis with multiplane images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 551–560, 2020. 2
- [46] Michael Waechter, Nils Moehrl, and Michael Goesele. Let there be color! large-scale texturing of 3d reconstructions.

- In *European conference on computer vision*, pages 836–850. Springer, 2014. [2](#)
- [47] Liao Wang, Jiakai Zhang, Xinhang Liu, Fuqiang Zhao, Yan-shun Zhang, Yingliang Zhang, Minye Wu, Jingyi Yu, and Lan Xu. Fourier plencotrees for dynamic radiance field rendering in real-time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13524–13534, 2022. [2](#), [8](#)
- [48] Daniel N Wood, Daniel I Azuma, Ken Aldinger, Brian Curless, Tom Duchamp, David H Salesin, and Werner Stuetzle. Surface light fields for 3d photography. In *Proceedings of the 27th annual conference on Computer graphics and interactive techniques*, pages 287–296, 2000. [2](#)
- [49] 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*, pages 9421–9431, 2021. [1](#), [2](#)
- [50] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *arXiv preprint arXiv:2112.05131*, 2021. [2](#), [3](#), [6](#)
- [51] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plencotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. [2](#), [3](#)
- [52] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. [2](#)
- [53] 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*, pages 586–595, 2018. [6](#)
- [54] 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. [2](#)
- [55] C Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski. High-quality video view interpolation using a layered representation. *ACM transactions on graphics (TOG)*, 23(3):600–608, 2004. [1](#), [2](#)