

4K4D: Real-Time 4D View Synthesis at 4K Resolution

Zhen Xu¹ Sida Peng¹ Haotong Lin¹ Guangzhao He¹
Jiaming Sun¹ Yujun Shen² Hujun Bao¹ Xiaowei Zhou^{1*}
¹Zhejiang University ²Ant Group



Figure 1. **Photorealistic and real-time rendering of dynamic 3D scenes.** Our proposed method reconstructs a 4D neural representation from multi-view videos, which can be rendered at 1125×1536 resolution with a speed of over 200 FPS using an RTX 3090 GPU while maintaining state-of-the-art quality on the DNA-Rendering [12] dataset. It is also noteworthy that our method reaches over 80 FPS when rendering 4K images with an RTX 4090. Detailed performance under different resolutions using different GPUs can be found in Tab. 5.

Abstract

This paper targets high-fidelity and real-time view synthesis of dynamic 3D scenes at 4K resolution. Recent methods on dynamic view synthesis have shown impressive rendering quality. However, their speed is still limited when rendering high-resolution images. To overcome this problem, we propose 4K4D, a 4D point cloud representation that supports hardware rasterization and network pre-computation to enable unprecedented rendering speed with a high rendering quality. Our representation is built on a 4D feature grid so that the points are naturally regularized and can be robustly optimized. In addition, we design a novel hybrid appearance model that significantly boosts the rendering quality while preserving efficiency. Moreover, we develop a differentiable depth peeling algorithm to effectively learn the proposed model from RGB videos. Experiments show that our representation can be rendered at over 400 FPS on the DNA-Rendering dataset at 1080p resolution and 80 FPS on the ENeRF-Outdoor dataset at 4K resolution using an RTX 4090 GPU, which is $30\times$ faster than previous methods and achieves the state-of-the-art rendering quality. Our project page is available at <https://zju3dv.github.io/4k4d>.

The authors from Zhejiang University are affiliated with the State Key Lab of CAD&CG. *Corresponding author: Xiaowei Zhou.

1. Introduction

Dynamic view synthesis aims to reconstruct dynamic 3D scenes from captured videos and create free-viewpoint and immersive virtual playback, which is a long-standing research problem in computer vision and computer graphics. Essential to the practicality of this technique is its ability to be rendered in real-time with high fidelity. Traditional methods [7, 13, 15, 16, 26, 61, 62, 84, 99, 100] represent dynamic 3D scenes as textured mesh sequences which can be rendered efficiently. However, high-quality mesh reconstruction requires complicated capture hardware and is limited to controlled environments.

Recently, implicit neural representations [19, 42, 58] have shown great success in reconstructing dynamic 3D scenes from RGB videos via differentiable rendering. For example, Li *et al.* [42] model the target scene as a dynamic neural radiance field and leverage volume rendering [17] to synthesize images. Despite impressive view synthesis results, existing approaches typically require seconds or even minutes to render an image at 1080p resolution due to the costly network evaluation, as discussed by Peng *et al.* [68]. Inspired by static view synthesis approaches [20, 33, 97], some dynamic view synthesis methods [2, 49, 68, 89] increase the rendering speed by decreasing either the

network size or the number of network evaluations. With these strategies, such methods achieve over 40 FPS when rendering moderate-resolution images (384×512) [49, 68], but are still not fast enough to achieve real-time performance when rendering high-resolution images. For instance, when rendering 4K resolution images, their speed reduces to only 1 or 2 FPS [2, 49, 68].

In this paper, we propose a novel neural representation, named 4K4D, for modeling and rendering dynamic 3D scenes. As illustrated in Fig. 1, 4K4D significantly outperforms previous dynamic view synthesis approaches [19, 49] in terms of the rendering speed, while being competitive in the rendering quality. Our core innovation lies in a 4D point cloud representation and a hybrid appearance model.

Specifically, for the dynamic scene, we obtain the coarse point cloud sequence using space carving [37] and model the position of each point as a learnable vector. A 4D feature grid is introduced for assigning a feature vector to each point, which is fed into MLP networks to predict the point’s radius, density, and spherical harmonics (SH) coefficients [59]. The 4D feature grid naturally applies spatial regularization on the point clouds and makes the optimization more robust (Sec. 5.2). During inference, the point’s radius, density and SH coefficients can be pre-computed, which eliminates network evaluations to achieve unprecedented rendering speed. Moreover, we develop a differentiable depth peeling algorithm that exploits the hardware rasterizer to further significantly accelerate the rendering.

We empirically find that the image blending model [49] achieves higher rendering quality than the SH model used by 3DGS [33]. However, the image blending model of previous methods [48, 49, 90] requires slow network evaluations during inference, limiting their rendering speed. To alleviate this, we introduce a novel design where we make the image blending network independent of the viewing direction, so the network evaluation can be pre-computed and thereby boost the rendering speed. As a two-edged sword, this strategy makes the appearance model discrete along the viewing direction. This downside is compensated for by using another continuous SH model.

To validate the effectiveness of the proposed pipeline, we evaluate 4K4D on multiple widely used datasets for multi-view dynamic novel view synthesis, including NHR [93], ENeRF-Outdoor [49], DNA-Rendering [12], and Neural3DV [41]. Extensive experiments show that 4K4D could not only be rendered orders of magnitude faster but also notably outperform the baselines in terms of rendering quality. With an RTX 4090 GPU, our method reaches 400 FPS on the DNA-Rendering dataset at 1080p resolution and 80 FPS on the ENeRF-Outdoor dataset at 4K resolution.

2. Related Work

Traditional scene representations. In the domain of novel view synthesis, various approaches based on different representations have been proposed, including multi-view image-based methods [6, 8, 18, 31, 69, 103], multi-plane image representations [47, 56, 65, 83, 86, 86], light-field techniques [14, 21, 39] as well as explicit surface or voxel-based methods [5, 13, 15, 22, 44, 61, 62, 100]. The seminal work [13] utilizes depth sensors and multi-view stereo techniques to consolidate per-view depth information into a coherent mesh sequence, producing high-quality volumetric video. These methods require intricate hardware setups and studio arrangements, thus constraining their accessibility.

Neural scene representations. Recently, implicit neural scene representations [3, 24, 27, 30, 32, 51, 52, 58, 76, 79–81, 85, 91] have attracted significant interest among researchers. NeRF [58] encodes the radiance fields of static scenes using coordinate-based Multi-Layer Perceptrons (MLP), achieving exceptional novel view synthesis quality. Building upon NeRF, a collection of studies [28, 42, 45, 63, 64, 70, 93] have made extensions to accommodate for dynamic scenes. Another line of studies [10, 46, 90, 98] has focused on integrating image features into the NeRF rendering pipeline. This approach is easily applicable to dynamic scenes, as multi-view videos can be directly decomposed into multi-view images. However, NeRF-based approaches often suffer from substantial network evaluation costs during the volume rendering process, which significantly limits their rendering speed and thus hinders their practicality.

Accelerating neural scene representations. To accelerate NeRF’s rendering, multiple works propose to distill implicit MLP networks into explicit structures that offer fast query capabilities, including voxel grids [20, 25, 40, 60, 72, 96, 97], explicit surfaces [11, 23, 29, 36, 54, 67] and point-based representations [1, 33, 35, 38, 71, 73, 101]. These methods effectively reduce the cost or the number of NeRF’s MLP evaluations required. Inspired by their success, several approaches [2, 9, 48, 49, 53, 68, 75, 82, 82, 87, 88] have explored the possibility of real-time dynamic view synthesis. HyperReel [2] employs a primitive prediction module to reduce the number of network evaluations, thereby achieving real-time speed at moderate resolutions. However, it should be noted that their rendering speed decreases significantly when rendering higher-resolution images.

Gaussian Splatting. One notable advancement for accelerating NeRF is the development of 3D Gaussian Splatting (3DGS) [33] which introduces a differentiable Gaussian ellipsoids splatting algorithm for fast and differentiable volume rendering [4, 17]. By effectively eliminating the slow ray marching operation of NeRF with forward splatting and SH [59], they attain both high-fidelity and high-speed rendering. However, the storage cost of 3DGS limits its application

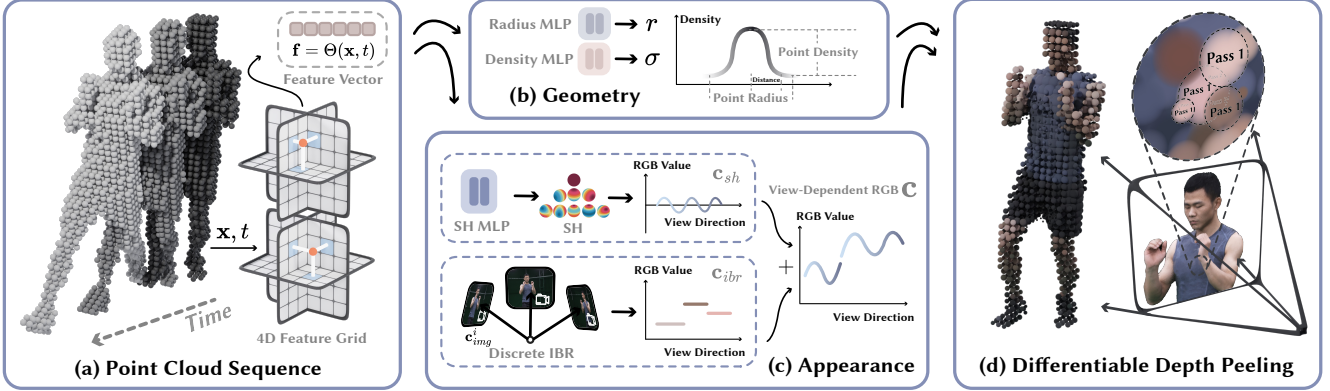


Figure 2. **Overview of our proposed pipeline.** (a) By applying the space-carving algorithm [37], we extract the initial cloud sequence \mathbf{x}, t of the target scene. A 4D feature grid [19] is predefined to assign a feature vector to each point, which is then fed into MLPs for the scene geometry and appearance. (b) The geometry model is based on the point location, radius, and density, which forms a semi-transparent point cloud. (c) The appearance model consists of a piece-wise constant IBR term \mathbf{c}_{ibr} and a continuous SH model \mathbf{c}_{sh} . (d) The proposed representation is learned from multi-view RGB videos through the differentiable depth peeling algorithm.

on dynamic scenes. In contrast, the 4D feature grid and image blending model of 4K4D could not only maintain similar rendering quality but also significantly reduce the storage cost for modeling dynamic scenes. Moreover, the simpler point cloud representation and the 4D feature grid regularization also make 4K4D less prone to overfitting training views than 3DGS. Some recent concurrent works [43, 55, 92, 94, 95] have also reported real-time rendering speeds by incorporating temporal correspondence or time-dependency into 3DGS. However, these methods either do not show results on datasets with large and fast motions [43, 95] (like NHR [93]) or only report real-time speed at moderate resolution (800×800 [92] and 640×480 [55]). In contrast, 4K4D is capable of real-time rendering even at 4K resolution while concurrently maintaining state-of-the-art view-synthesis quality on large-motion data.

3. Proposed Approach

Given a multi-view video capturing a dynamic 3D scene, our goal is to reconstruct the target scene and perform novel view synthesis in real time. To this end, we extract coarse point clouds of the scene using the space-carving algorithm [37] (Sec. 4) and build a point cloud-based neural scene representation, which can be robustly learned from input videos and enable the hardware-accelerated rendering.

The overview of the proposed model is presented in Fig. 2. In this section, we first describe how to represent the geometry and appearance of dynamic scenes based on point clouds and neural networks (Sec. 3.1). Then, we develop a differentiable depth peeling algorithm for rendering our representation (Sec. 3.2), which is supported by the hardware rasterizer, thereby significantly improving the rendering speed. Finally, we discuss how to optimize the proposed model on input RGB videos (Sec. 3.3).

3.1. Modeling Dynamic Scenes with Point Clouds

4D embedding. Given the coarse point clouds of the target scene, we represent its dynamic geometry and appearance using neural networks and feature grids. Specifically, our method first defines six feature planes $\theta_{xy}, \theta_{xz}, \theta_{yz}, \theta_{tx}, \theta_{ty}$, and θ_{tz} . To assign a feature vector \mathbf{f} to any point \mathbf{x} at frame t , we adopt the strategy of K-Planes [19] to model a 4D feature field $\Theta(\mathbf{x}, t)$ using these six planes:

$$\mathbf{f} = \Theta(\mathbf{x}, t) = \theta_{xy}(x, y) \oplus \theta_{xz}(x, z) \oplus \theta_{yz}(y, z) \oplus \theta_{tx}(t, x) \oplus \theta_{ty}(t, y) \oplus \theta_{tz}(t, z), \quad (1)$$

where $\mathbf{x} = (x, y, z)$ is the input point, and \oplus indicates the concatenation operator. Please refer to K-Planes [19] for more implementation details.

Geometry model. Based on coarse point clouds, the dynamic scene geometry is represented by learning three entries on each point: position $\mathbf{p} \in R^3$, radius $r \in R$, and density $\sigma \in R$. Using these point entries, we calculate the volume density of space point \mathbf{x} with respect to an image pixel \mathbf{u} for the volume rendering, which will be described in Sec. 3.2. The point position \mathbf{p} is modeled as an optimizable vector. The radius r and density σ are predicted by feeding the feature vector \mathbf{f} in Eq. (1) to an MLP network.

Appearance model. As illustrated in Fig. 2c, we use the image blending technique and the spherical harmonics (SH) model [59, 97] to build a hybrid appearance model, where the image blending technique represents the discrete view-dependent appearance \mathbf{c}_{ibr} and the SH model represents the continuous view-dependent appearance \mathbf{c}_{sh} . For point \mathbf{x} at frame t , its color with viewing direction \mathbf{d} is:

$$\mathbf{c}(\mathbf{x}, t, \mathbf{d}) = \mathbf{c}_{ibr}(\mathbf{x}, t, \mathbf{d}) + \mathbf{c}_{sh}(\mathbf{s}, \mathbf{d}), \quad (2)$$

where \mathbf{s} means SH coefficients at point \mathbf{x} .

The discrete view-dependent appearance \mathbf{c}_{ibr} is inferred based on input images. Specifically, for a point \mathbf{x} , we first project it into the input image to retrieve the corresponding RGB color \mathbf{c}_{img}^i . Then, to blend input RGB colors, we calculate the corresponding blending weight w^i based on the point coordinate and the input image. Note that the blending weight is independent from the viewing direction. Next, to achieve the view-dependent effect, we select the N' nearest input views according to the viewing direction. Finally, the color \mathbf{c}_{ibr} is computed as $\sum_{i=1}^{N'} w^i \mathbf{c}_{img}^i$. Because the N' input views are obtained through the nearest neighbor retrieval, the \mathbf{c}_{ibr} is inevitably discrete along the viewing direction. To achieve the continuous view-dependent effect, we append the fine-level color \mathbf{c}_{sh} represented by the SH model, as shown in Fig. 2c.

In practice, our method regresses the SH coefficients \mathbf{s} by passing the point feature \mathbf{f} in Eq. (1) into an MLP network. To predict the blending weight w^i in the image blending model \mathbf{c}_{ibr} , we first project point \mathbf{x} onto the input image to retrieve the image feature \mathbf{f}_{img}^i , and then concatenate it with the point feature \mathbf{f} , which is fed into another MLP network to predict the blending weight. The image feature \mathbf{f}_{img}^i is extracted using a 2D CNN network.

Discussion. Our appearance model is the key to achieving the low-storage, high-fidelity, and real-time view synthesis of dynamic scenes. There are three alternative ways to represent the dynamic appearance, but they cannot perform on par with our model. 1) Defining explicit SH coefficients on each point, as in 3D Gaussian splatting [33]. When the degree of SH coefficients is high and the amount of points of dynamic scenes is large, this model’s size could be too big to train on a consumer GPU. 2) MLP-based SH model. Using an MLP to predict SH coefficients of each point can effectively decrease the model size. However, our experiments found that MLP-based SH model struggles to render high-quality images (Sec. 5.2). 3) Continuous view-dependent image blending model, as in ENeRF [49]. We found that representing the appearance with the image blending model exhibits better rendering quality than only with the MLP-based SH model. However, the color network in ENeRF takes the viewing direction as input and thus cannot be easily pre-computed, limiting the rendering speed during inference.

In contrast to these three methods, our appearance model combines a discrete image blending model \mathbf{c}_{ibr} with a continuous SH model \mathbf{c}_{sh} . The image blending model \mathbf{c}_{ibr} boosts the rendering performance. In addition, it supports the pre-computation, as its network does not take the viewing direction as input. The SH model \mathbf{c}_{sh} enables the view-dependent effect for any viewing direction. During training, our model represents the scene appearance using networks, so its model size is reasonable. During inference, we pre-compute the network outputs to achieve the real-time rendering, which will be described in Sec. 3.4.

3.2. Differentiable Depth Peeling

Our proposed dynamic scene representation can be rendered into images by performing volume rendering [17] on rasterized points. This forward process is much faster than NeRF’s backward ray-marching operation [57] since it requires no network evaluation and explicit sampling. The volume rendering equation requires the color and transparency values to be integrated in order [4], thus we utilize the depth-peeling algorithm for acquiring the corresponding ordered points for pixels. Thanks to the point cloud representation, we can leverage the hardware rasterizer to significantly speed up the depth peeling and blending process. Moreover, it is easy to make this rendering process differentiable, enabling us to learn our model from input RGB videos.

We develop a custom shader to implement the depth peeling algorithm that consists of K rendering passes. Consider a particular image pixel \mathbf{u} . In the first pass, our method first uses the hardware rasterizer to render point clouds onto the image, which assigns the closest-to-camera point \mathbf{x}_0 to the pixel \mathbf{u} . Denote the depth of point \mathbf{x}_0 as t_0 . Subsequently, in the k -th rendering pass, all points with depth value t_k smaller than the recorded depth of the previous pass t_{k-1} are discarded, thereby resulting in the k -th closest-to-camera point \mathbf{x}_k for the pixel \mathbf{u} . Discarding closer points is implemented in our custom shader, so it still supports the hardware rasterization. After K rendering passes, pixel \mathbf{u} has a set of sorted points $\{\mathbf{x}_k | k = 1, \dots, K\}$.

Based on the sorted points, we use the volume rendering technique to synthesize the color of pixel \mathbf{u} . The densities of these points for pixel \mathbf{u} are defined based on the distance between the projected point and pixel \mathbf{u} on the 2D image:

$$\alpha(\mathbf{u}, \mathbf{x}) = \sigma \cdot \max\left(1 - \frac{\|\pi(\mathbf{x}) - \mathbf{u}\|_2^2}{r^2}, 0\right), \quad (3)$$

where π is the camera projection function. σ and r are the density and radius of point \mathbf{x} , which are described in Sec. 3.1. Intuitively, Eq. (3) defines a semi-transparent point representation where the density is the highest around the center and quadratically decreases along its radius. During training, we implement the projection function π using the PyTorch [66], so Eq. (3) is naturally differentiable. During inference, we leverage the hardware rasterization process to efficiently obtain the distance $\|\pi(\mathbf{x}) - \mathbf{u}\|_2^2$, which is implemented using OpenGL [77].

Denote the density of point \mathbf{x}_k as α_k . The color of pixel \mathbf{u} from the volume rendering is formulated as:

$$C(\mathbf{u}) = \sum_{k=1}^K T_k \alpha_k \mathbf{c}_k, \text{ where } T_k = \prod_{j=1}^{k-1} (1 - \alpha_j), \quad (4)$$

where \mathbf{c}_k is the color of point \mathbf{x}_k , as described in Eq. (2).

3.3. Training

Given the rendered pixel color $C(\mathbf{u})$, we compare it with the ground-truth pixel color $C_{gt}(\mathbf{u})$ to optimize our model in an end-to-end fashion using the following loss function:

$$L_{img} = \sum_{\mathbf{u} \in \mathcal{U}} \|C(\mathbf{u}) - C_{gt}(\mathbf{u})\|_2^2, \quad (5)$$

where \mathcal{U} is the set of image pixels. In addition to the MSE loss L_{img} , we also apply the perceptual loss $L_{lpi ps}$ [102].

$$L_{lpi ps} = \|\Phi(I) - \Phi(I_{gt})\|_1, \quad (6)$$

where Φ is the perceptual function (a VGG16 network) and I, I_{gt} are the rendered and ground-truth images, respectively. The perceptual loss [102] computes the difference in image features extracted from the VGG model [78]. Our experiments in Sec. 5.2 show that it effectively improves the perceived quality of the rendered image.

To regularize the optimization process of our proposed representation, we additionally apply mask supervision to dynamic regions of the target scene. We solely render point clouds of dynamic regions to obtain their masks, where the pixel value is obtained by:

$$M(\mathbf{u}) = \sum_{k=1}^K T_k \alpha_k, \text{ where } T_k = \prod_{j=1}^{k-1} (1 - \alpha_j). \quad (7)$$

The mask loss is defined as:

$$L_{msk} = - \sum_{\mathbf{u} \in \mathcal{U}'} \|M(\mathbf{u}) - M_{gt}(\mathbf{u})\|_2^2, \quad (8)$$

where \mathcal{U}' means the set of pixels of the rendered mask, and M_{gt} is the ground-truth mask of 2D dynamic regions. This effectively regularizes the optimization of the geometry of dynamic regions by confining it to the visual hulls.

The final loss function is defined as

$$L = L_{img} + \lambda_{lpi ps} L_{lpi ps} + \lambda_{msk} L_{msk}, \quad (9)$$

where $\lambda_{lpi ps}$ and λ_{msk} are hyperparameters controlling weights of correspondings losses.

3.4. Inference

After training, we apply a few acceleration techniques to boost the rendering speed of our model. First, we precompute the point location \mathbf{p} , radius r , density σ , SH coefficients s and color blending weights w_i before inference, which are stored at the main memory. During rendering, these properties are asynchronously streamed onto the graphics card, overlapping rasterization with memory copy to achieve an optimal rendering speed [74, 77]. After applying this technique, the runtime computation is reduced to only a

depth peeling evaluation (Sec. 3.2) and a spherical harmonics evaluation (Eq. (2)). Second, we convert the model from 32-bit floats to 16-bits for efficient memory access, which increases FPS by 20 and leads to no visible performance loss. Third, the number of rendering passes K for the differentiable depth peeling algorithm is reduced from 15 to 12, also leading to a 20 FPS speedup with no visual quality change. Detailed analyses of rendering speed can be found in Sec. 5.2 and the supplementary material.

4. Implementation Details

Optimization. 4K4D is trained using the PyTorch framework [66]. Using the Adam optimizer [34] with a learning rate $5e^{-3}$, our models typically converge after 800k iterations for a sequence length of 200 frames, which takes around 24 hours on a single RTX 4090 GPU. Specifically, the learning rate of point positions is set to $1e^{-5}$, and the regularization loss weights $\lambda_{lpi ps}$ and λ_{msk} are set to $1e^{-3}$. During training, the number of passes K for the differentiable depth peeling is set to 15, and the number of nearest input views N' is set to 4. The rendering speed of our method is reported on an RTX 3090 GPU for the experiments in Sec. 5 unless otherwise stated.

Initialization of point clouds. We leverage existing multi-view reconstruction methods to initialize the point clouds. For dynamic regions, we use segmentation methods [50] to obtain their masks in input images and utilize the space carving algorithm [37] to extract their coarse geometry. For static background regions, we leverage foreground masks to compute the mask-weighted average of background pixels along all frames, producing background images without the foreground content. Then, an Instant-NGP [60] model is trained on these images, from which we obtain the initial point clouds. After the initialization, the number of points for the dynamic regions is typically 250k per frame, and the static background regions typically consist of 300k points.

5. Experiments

Datasets. We train and evaluate our method 4K4D on multiple widely used multi-view datasets, including DNA-Rendering [12], ENeRF-Outdoor [49] and NHR [93]. DNA-Rendering [12] records 10-second clips of dynamic humans and objects at 15 FPS using 4K and 2K cameras with 60 views. This dataset is very challenging due to the complex clothing and fast motions. We conduct experiments on 4 sequences of DNA-Rendering, with 90% of the views as training set and the rest as evaluation set. ENeRF-Outdoor [49] records multiple dynamic humans and objects in an outdoor environment at 30FPS using 1080p cameras. We select three 100-frame sequences with 6 different actors (2 for each sequence) holding objects for evaluation. This dataset is difficult for dynamic view synthesis in that not

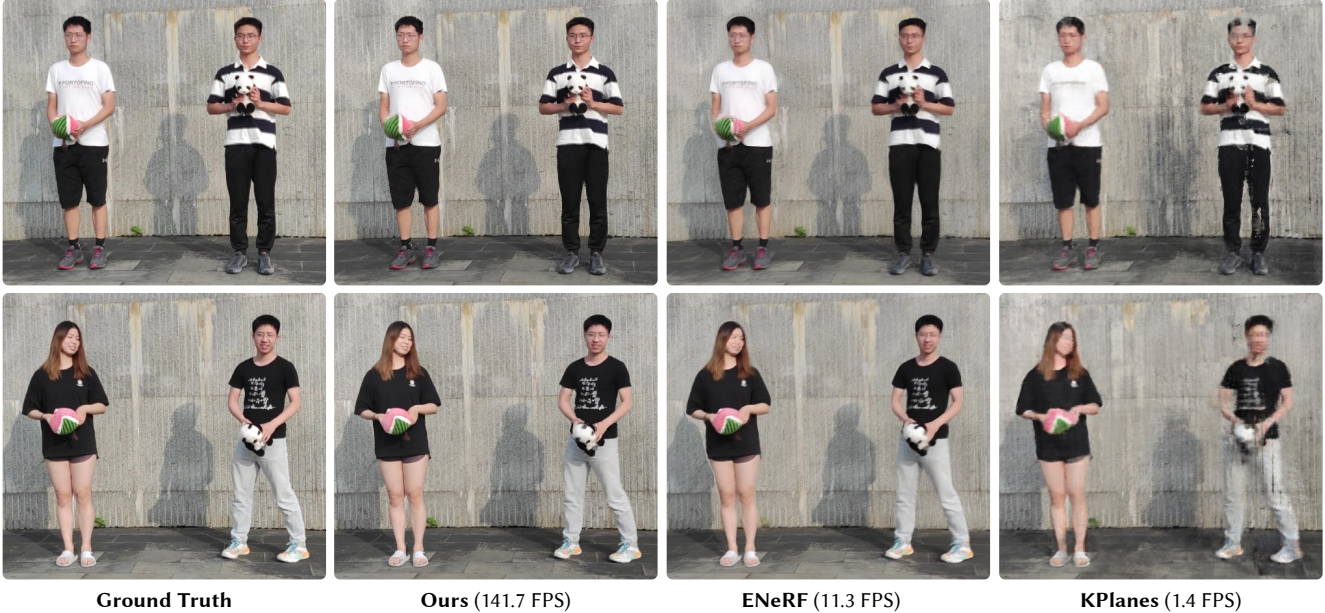


Figure 3. **Qualitative comparison on the ENeRF-Outdoor [49] dataset that contains 960×540 images.** Our method achieves much higher rendering quality and can be rendered $14\times$ faster than ENeRF[49]. More dynamic results can be found in the supplementary video.

only are there multiple moving humans and objects, but the background is also dynamic due to cast shadows. More details can be found in the supplementary.

5.1. Comparison Experiments

Comparison on DNA-Rendering [12]. Qualitative and quantitative comparisons on DNA-Rendering [12] are shown in Fig. 4 and Tabs. 1 and 3 respectively. As evident in Tab. 1, our method renders $30\times$ faster than the SOTA real-time dynamic view synthesis method ENeRF [49] with superior quality. Even when compared with concurrent work [48], our method still achieves $13\times$ speedup and produces consistently higher quality images. As shown in Fig. 4, KPlanes [19] could not recover the highly detailed appearance and geometry of the 4D dynamic scene. Other image-based methods [48, 49, 90] produce high-quality appearance. However, they tend to produce blurry results around occlusions and edges, leading to degradation of the visual quality while maintaining interactive framerate at best. When compared with 3DGS [33] on the first frame of each sequence, our method achieves a much better storage efficiency ($50\times$) thanks to our compact 4D feature grid and image blending model. Moreover, due to the simplicity of our point-based representation, our method is less prone to overfit the training views. More details of the comparison with 3DGS can be found in the supplementary material.

Comparison on ENeRF-Outdoor [49]. Fig. 3 and Tabs. 2 and 3 provides qualitative and quantitative results on the ENeRF-Outdoor [49] dataset. Even on the challenging ENeRF-Outdoor dataset with multiple actors and the back-

Table 1. **Quantitative comparison on the DNA-Rendering [12] dataset.** Image resolutions are 1024×1224 and 1125×1536 . Metrics are averaged over all scenes. Green and yellow cell colors indicate the best and the second best results, respectively.

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FPS
ENeRF [49]	28.108	0.972	0.056	6.011
IBRNet [90]	27.844	0.967	0.081	0.100
KPlanes [19]	27.452	0.952	0.118	0.640
Im4D [48]	28.991	0.973	0.062	15.360
Ours	31.173	0.976	0.055	203.610

Table 2. **Quantitative comparison on the ENeRF-Outdoor [49] dataset.** This dataset includes 960×540 images. Green and yellow cell colors indicate the best and the second-best results, respectively.

	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	FPS
ENeRF [49]	25.452	0.809	0.273	11.309
IBRNet [90]	24.966	0.929	0.172	0.140
KPlanes [19]	21.310	0.735	0.454	1.370
Ours	25.815	0.898	0.147	141.665

ground, our method still achieves notably better results while rendering at over 140 FPS. ENeRF [49] produces blurry results on this challenging dataset, and the rendering results of IBRNet [90] contain black artifacts around the edges of the images as shown in Fig. 3. K-Plane [19] fails to reconstruct the dynamic humans and varying background regions. 3DGS [33] not only introduces much higher storage cost than our method ($45\times$), but also faces even more pronounced overfitting problem with smaller number of views (18 for ENeRF-Outdoor). As evident in Tab. 3 and the

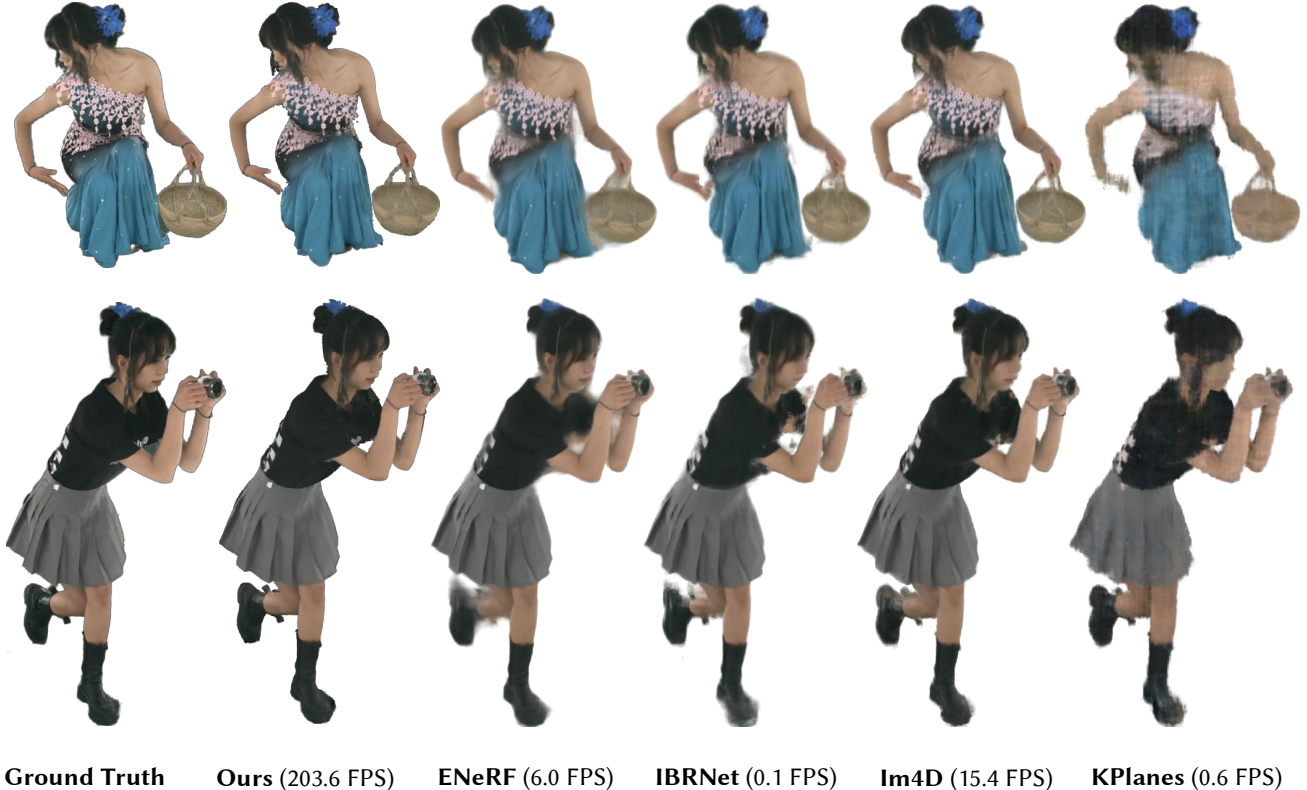


Figure 4. **Qualitative comparison on the DNA-Rendering [12] dataset that contains 1024×1224 (and 1125×1536) images.** Our method can produce high-fidelity images at over 200 FPS while other competitors fail to produce high-quality results for highly dynamic scenes.

supplementary material, the overfitting severely degrades the rendering quality. Their rendering speed is slower than ours due to excessive point count. More details of the comparison with 3DGS are present in the supplementary material.

5.2. Ablation Studies

We perform ablation studies on the proposed components on the 150-frame *0013_01* sequence of the DNA-Rendering [12] dataset. Our method can be rendered at over 200 FPS with state-of-the-art quality and maintains a only 2MB per frame storage overhead. More detailed rendering speed analysis and breakdown and storage cost analysis can be found in the supplementary material.

Ablation study on the 4D embedding. The “w/o f” variant removes the proposed 4D embedding (Sec. 3.1) module and replaces it with a per-frame and per-point optimizable position, radius, density, and scale. As shown in Fig. 5 and Tab. 4, the “w/o f” variant produces blurry and noisy geometry without the 4D embedding Θ , which leads to the inferior rendering quality.

Ablation study on the hybrid appearance model. The “w/o c_{ibr} ” variant removes c_{ibr} in the appearance formulation Eq. (2), which not only leads to less details on the recovered appearance but also significantly impedes the quality of the geometry. Adding an additional degree for the SH

Table 3. **Quantitative comparison on the first frame of all sequences of DNA-Rendering [12] (1024×1224 (and 1125×1536) images) and ENeRF-Outdoor [49] (960×540 images).** Metrics are averaged for each dataset. “Storage” indicates the disk file size of the trained models (including source images for our method).

Dataset	Method	PSNR	LPIPS	FPS	Storage	Training
DNA-Rendering	3DGS [33]	31.16	0.049	113.2	224 MB	5min
	Ours	31.87	0.046	241.7	4.7 MB	15min
ENeRF-Outdoor	3DGS [33]	21.63	0.349	88.4	715 MB	10min
	Ours	26.54	0.145	148.6	16.0 MB	30min

coefficients does not lead to a significant performance change (PSNR 30.129 vs. 30.259). Comparatively, our proposed method produces high-fidelity rendering with better details. A visualization of the view-dependent effect produced by c_{sh} can be found in the supplementary material.

Ablation study on loss functions. As shown in Tab. 4, removing the L_{lips} term not only reduces the perceptual quality (LPIPS score) but also leads to the degradation of other performance metrics. For the highly dynamic DNA-Rendering [12] dataset, the mask loss L_{msk} helps with regularizing the optimization of the dynamic geometry.

Rendering speed on different GPUs and resolutions. We additionally report the rendering speed of our method on different hardware (RTX 3060, RTX 3090, and RTX 4090)

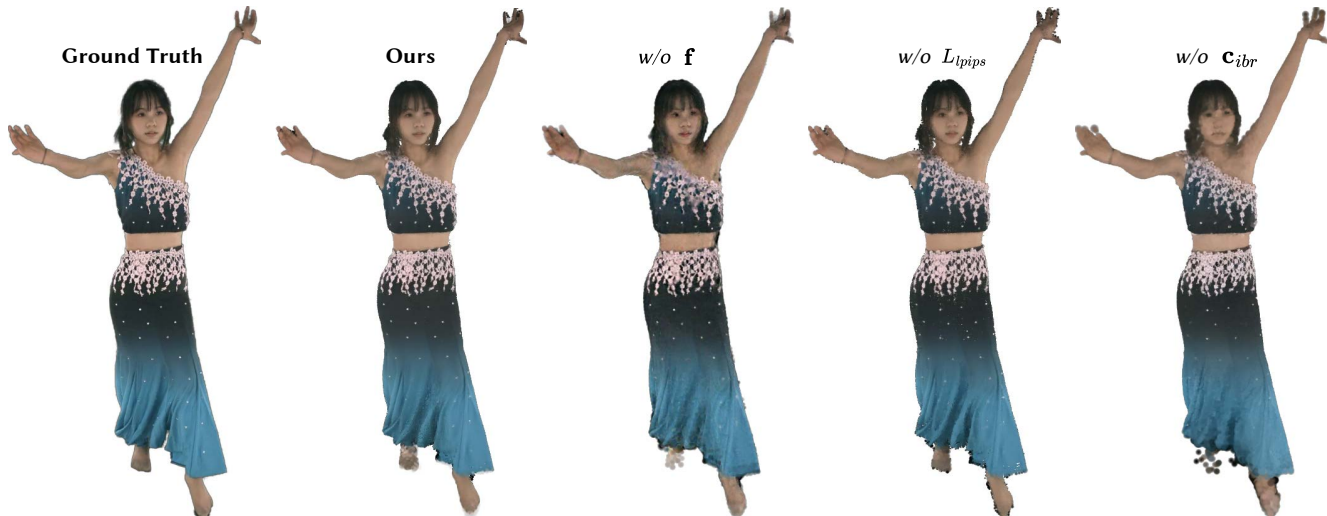


Figure 5. **Ablation studies** on the *0013_01* sequence of DNA-Rendering [12]. Removing our proposed components leads to noisy geometry and blurry appearance. Our method produces high-fidelity results with perceptually accurate shapes and colors. See Sec. 5.2 for more details.

with different resolutions (720p, 1080p, and 4K (2160p)) in Tab. 5. The rendering speed reported here contains the overhead of the interactive GUI. 4K4D achieves real-time rendering speed even when rendering 4K (2160p) images on commodity hardware as shown in the table. More real-time rendering demos can be found in the supplementary video.

6. Conclusion and Discussion

In this paper, we provide a neural point cloud-based representation, 4K4D, for real-time rendering of dynamic 3D scenes at 4K resolution. We build 4K4D upon a 4D feature grid to naturally regularize the points and develop a novel hybrid appearance model for high-quality rendering. Furthermore, we develop a differentiable depth peeling algorithm that utilizes the hardware rasterization pipeline to effectively optimize and efficiently render the proposed model. In our experiments, we demonstrate that 4K4D not only achieves state-of-the-art rendering quality but also exhibits a more than 30× increase in rendering speed (over 200FPS at 1080p on an RTX 3090 GPU).

However, our method still has some limitations. For one, 4K4D cannot produce correspondences of points across frames, which are important for some downstream tasks. Moreover, the storage cost for 4K4D increases linearly with the number of video frames, so our method has difficulty in modeling long volumetric videos. How to model correspondences and reduce the storage cost for long videos could be two interesting problems for future works. Moreover, the rendering quality of our method also depends on the resolution of input images. While our method achieves real-time rendering at 4K resolution, 4K-quality rendering can only be achieved with sufficient input resolution.

Table 4. **Ablation studies** on the 150-frame *0013_01* sequence of the DNA-Rendering dataset [12]. “w/o f” indicates replacing the 4D embedding with a per-frame and per-point optimizable position, radius, density, and scale. See Sec. 5.2 for more detailed descriptions for the abbreviations.

	PSNR ↑	SSIM ↑	LPIPS ↓	Model Size
w/o f	29.779	0.967	0.057	1304.0 MiB
w/o c_ibr	30.259	0.973	0.054	225.0 MiB
w/o c_sh	31.946	0.981	0.040	225.0 MiB
w/o L_lips	31.661	0.979	0.063	225.0 MiB
w/o L_msk	29.115	0.965	0.073	225.0 MiB
Ours	31.990	0.982	0.040	225.0 MiB

Table 5. **Rendering speed on different GPUs and resolutions.** The results are recorded on the first frame of the *0013_01* sequence of DNA-Rendering [12] and the *actor1_4* sequence of ENeRF-Outdoor [49] with the interactive GUI. Resolutions are set to 720p (720 × 1280), 1080p (1080 × 1920), and 4K (2160 × 3840). Even with the overhead of the interactive GUI (“w/ GUI”), our method still achieves unprecedented rendering speed. More real-time rendering results can be found in the supplementary video.

Dataset	Res.	RTX 3060	RTX 3090	RTX 4090
DNA-Rendering [12]	720p	173.8 FPS	246.9 FPS	431.0 FPS
	1080p	138.7 FPS	233.1 FPS	409.8 FPS
	4K	90.0 FPS	147.4 FPS	288.8 FPS
ENeRF-Outdoor [49]	720p	90.5 FPS	130.5 FPS	351.5 FPS
	1080p	66.1 FPS	103.6 FPS	249.7 FPS
	4K	25.1 FPS	47.2 FPS	85.1 FPS

Acknowledgement

The authors would like to acknowledge support from NSFC (No. 62172364) and Information Technology Center and State Key Lab of CAD&CG, Zhejiang University.

References

- [1] Kara-Ali Aliev, Artem Sevastopolsky, Maria Kolos, Dmitry Ulyanov, and Victor Lempitsky. Neural point-based graphics. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXII 16*, pages 696–712. Springer, 2020. [2](#)
- [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. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16610–16620, 2023. [1](#), [2](#)
- [3] Benjamin Attal, Jia-Bin Huang, Michael Zollhöfer, Johannes Kopf, and Changil Kim. Learning neural light fields with ray-space embedding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 19819–19829, 2022. [2](#)
- [4] Louis Bavoil and Kevin Myers. Order independent transparency with dual depth peeling. *NVIDIA OpenGL SDK*, 1:12, 2008. [2](#), [4](#)
- [5] Michael Broxton, John Flynn, Ryan Overbeck, Daniel Erickson, Peter Hedman, Matthew Duvall, Jason Dourgarian, Jay Busch, Matt Whalen, and Paul Debevec. Immersive light field video with a layered mesh representation. *ACM Transactions on Graphics (TOG)*, 39(4):86–1, 2020. [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, SIGGRAPH ’01*, page 425–432, New York, NY, USA, 2001. Association for Computing Machinery. [2](#)
- [7] Dan Casas, Marco Volino, John Collomosse, and Adrian Hilton. 4d video textures for interactive character appearance. In *Computer Graphics Forum*, pages 371–380. Wiley Online Library, 2014. [1](#)
- [8] Gaurav Chaurasia, Sylvain Duchene, Olga Sorkine-Hornung, and George Drettakis. Depth synthesis and local warps for plausible image-based navigation. *ACM TOG*, 2013. [2](#)
- [9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. *arXiv*, 2022. [2](#)
- [10] 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 *ICCV*, 2021. [2](#)
- [11] Zhiqin Chen, Thomas Funkhouser, Peter Hedman, and Andrea Tagliasacchi. Mobilenerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. *arXiv preprint arXiv:2208.00277*, 2022. [2](#)
- [12] Wei Cheng, Ruixiang Chen, Wanqi Yin, Siming Fan, Keyu Chen, Honglin He, Huiwen Luo, Zhongang Cai, Jingbo Wang, Yang Gao, et al. Dna-rendering: A diverse neural actor repository for high-fidelity human-centric rendering. *arXiv preprint arXiv:2307.10173*, 2023. [1](#), [2](#), [5](#), [6](#), [7](#), [8](#)
- [13] Alvaro Collet, Ming Chuang, Pat Sweeney, Don Gillett, Dennis Evseev, David Calabrese, Hugues Hoppe, Adam Kirk, and Steve Sullivan. High-quality streamable free-viewpoint video. *ACM Transactions on Graphics (ToG)*, 34(4):1–13, 2015. [1](#), [2](#)
- [14] Abe Davis, Marc Levoy, and Fredo Durand. Unstructured light fields. In *Computer Graphics Forum*, pages 305–314. Wiley Online Library, 2012. [2](#)
- [15] Mingsong Dou, Sameh Khamis, Yury Degtyarev, Philip Davidson, Sean Ryan Fanello, Adarsh Kowdle, Sergio Orts Escolano, Christoph Rhemann, David Kim, Jonathan Taylor, et al. Fusion4d: Real-time performance capture of challenging scenes. *ACM TOG*, 2016. [1](#), [2](#)
- [16] Mingsong Dou, Jonathan Taylor, Henry Fuchs, Andrew Fitzgibbon, and Shahram Izadi. 3d scanning deformable objects with a single rgbd sensor. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 493–501, 2015. [1](#)
- [17] Robert A Drebin, Loren Carpenter, and Pat Hanrahan. Volume rendering. *ACM Siggraph Computer Graphics*, 22(4):65–74, 1988. [1](#), [2](#), [4](#)
- [18] John Flynn, Ivan Neulander, James Philbin, and Noah Snavely. Deepstereo: Learning to predict new views from the world’s imagery. In *CVPR*, June 2016. [2](#)
- [19] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. [1](#), [2](#), [3](#), [6](#)
- [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*, pages 14346–14355, 2021. [1](#), [2](#)
- [21] Steven J Gortler, Radek Grzeszczuk, Richard Szeliski, and Michael F Cohen. The lumigraph. In *SIGGRAPH*, 1996. [2](#)
- [22] Kaiwen Guo, Peter Lincoln, Philip Davidson, Jay Busch, Xueming Yu, Matt Whalen, Geoff Harvey, Sergio Orts-Escolano, Rohit Pandey, Jason Dourgarian, et al. The relightables: Volumetric performance capture of humans with realistic relighting. *ACM Transactions on Graphics (ToG)*, 38(6):1–19, 2019. [2](#)
- [23] Jon Hasselgren, Nikolai Hofmann, and Jacob Munkberg. Shape, light, and material decomposition from images using monte carlo rendering and denoising. *Advances in Neural Information Processing Systems*, 35:22856–22869, 2022. [2](#)
- [24] Peter Hedman, Julien Philip, True Price, Jan-Michael Frahm, George Drettakis, and Gabriel Brostow. Deep blending for free-viewpoint image-based rendering. *ACM TOG*, 2018. [2](#)
- [25] Peter Hedman, Pratul P. Srinivasan, Ben Mildenhall, Jonathan T. Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. In *ICCV*, 2021. [2](#)
- [26] Anna Hilsmann, Philipp Fechteler, Wieland Morgenstern, Wolfgang Paier, Ingo Feldmann, Oliver Schreer, and Peter Eisert. Going beyond free viewpoint: creating animatable volumetric video of human performances. *IET Computer Vision*, pages 350–358, 2020. [1](#)
- [27] Tao Hu, Tao Yu, Zerong Zheng, He Zhang, Yebin Liu, and Matthias Zwicker. Hvtr: Hybrid volumetric-textural rendering for human avatars. In *2022 International Conference on 3D Vision (3DV)*, pages 197–208. IEEE, 2022.

- 2
- [28] Mustafa Işık, Martin Rüinz, Markos Georgopoulos, Taras Khakhulin, Jonathan Starck, Lourdes Agapito, and Matthias Nießner. Humanrf: High-fidelity neural radiance fields for humans in motion. *arXiv preprint arXiv:2305.06356*, 2023. 2
- [29] Shubhendu Jena, Franck Multon, and Adnane Boukhayma. Neural mesh-based graphics. In *European Conference on Computer Vision*, pages 739–757. Springer, 2022. 2
- [30] Yue Jiang, Dantong Ji, Zhizhong Han, and Matthias Zwicker. Sdfdiff: Differentiable rendering of signed distance fields for 3d shape optimization. In *CVPR*, 2020. 2
- [31] Nima Khademi Kalantari, Ting-Chun Wang, and Ravi Ramamoorthi. Learning-based view synthesis for light field cameras. *ACM TOG*, 2016. 2
- [32] Petr Kellnhofer, Lars C Jebe, Andrew Jones, Ryan Spicer, Kari Pulli, and Gordon Wetzstein. Neural lumigraph rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4287–4297, 2021. 2
- [33] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Transactions on Graphics (TOG)*, 42(4):1–14, 2023. 1, 2, 4, 6, 7
- [34] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [35] Georgios Kopanas, Julien Philip, Thomas Leimkühler, and George Drettakis. Point-based neural rendering with per-view optimization. In *Computer Graphics Forum*, volume 40, pages 29–43. Wiley Online Library, 2021. 2
- [36] Jonas Kulhanek and Torsten Sattler. Tetra-nerf: Representing neural radiance fields using tetrahedra. *arXiv preprint arXiv:2304.09987*, 2023. 2
- [37] Kiriakos N Kutulakos and Steven M Seitz. A theory of shape by space carving. *International journal of computer vision*, 38:199–218, 2000. 2, 3, 5
- [38] Christoph Lassner and Michael Zollhofer. Pulsar: Efficient sphere-based neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1440–1449, 2021. 2
- [39] Marc Levoy and Pat Hanrahan. Light field rendering. In *SIGGRAPH*, 1996. 2
- [40] Ruilong Li, Hang Gao, Matthew Tancik, and Angjoo Kanazawa. Nerfacc: Efficient sampling accelerates nerfs. *arXiv preprint arXiv:2305.04966*, 2023. 2
- [41] Tianye Li, Mira Slavcheva, Michael Zollhofer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, and Zhaoyang Lv. Neural 3d video synthesis. *arXiv preprint arXiv:2103.02597*, 2021. 2
- [42] Tianye Li, Mira Slavcheva, Michael Zollhofer, 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
- [43] Zhan Li, Zhang Chen, Zhong Li, and Yi Xu. Spacetime gaussian feature splatting for real-time dynamic view synthesis. *arXiv preprint arXiv:2312.16812*, 2023. 3
- [44] Zhong Li, Yu Ji, Wei Yang, Jinwei Ye, and Jingyi Yu. Robust 3d human motion reconstruction via dynamic template construction. In *2017 International Conference on 3D Vision (3DV)*, pages 496–505. IEEE, 2017. 2
- [45] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 2
- [46] Zhengqi Li, Qianqian Wang, Forrester Cole, Richard Tucker, and Noah Snavely. Dynibar: Neural dynamic image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4273–4284, 2023. 2
- [47] Zhengqi Li, Wenqi Xian, Abe Davis, and Noah Snavely. Crowdsampling the plenoptic function. In *ECCV*, 2020. 2
- [48] Haotong Lin, Sida Peng, Zhen Xu, Tao Xie, Xingyi He, Hujun Bao, and Xiaowei Zhou. High-fidelity and real-time novel view synthesis for dynamic scenes. In *SIGGRAPH Asia Conference Proceedings*, 2023. 2, 6
- [49] Haotong Lin, Sida Peng, Zhen Xu, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Efficient neural radiance fields for interactive free-viewpoint video. In *SIGGRAPH Asia Conference Proceedings*, 2022. 1, 2, 4, 5, 6, 7, 8
- [50] Shanchuan Lin, Linjie Yang, Imran Saleemi, and Soumyadip Sengupta. Robust high-resolution video matting with temporal guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 238–247, 2022. 5
- [51] Shichen Liu, Shunsuke Saito, Weikai Chen, and Hao Li. Learning to infer implicit surfaces without 3d supervision. *NeurIPS*, 2019. 2
- [52] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. In *SIGGRAPH*, 2019. 2
- [53] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhofer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering. *ACM Transactions on Graphics (TOG)*, 40(4):1–13, 2021. 2
- [54] Fan Lu, Yan Xu, Guang Chen, Hongsheng Li, Kwan-Yee Lin, and Changjun Jiang. Urban radiance field representation with deformable neural mesh primitives. *arXiv preprint arXiv:2307.10776*, 2023. 2
- [55] Jonathon Luiten, Georgios Kopanas, Bastian Leibe, and Deva Ramanan. Dynamic 3d gaussians: Tracking by persistent dynamic view synthesis. *arXiv preprint arXiv:2308.09713*, 2023. 3
- [56] 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 TOG*, 2019. 2
- [57] 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. 2020. 4
- [58] 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
- [59] Claus Müller. *Spherical harmonics*, volume 17. Springer, 2006. 2, 3
- [60] 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):1–15, 2022. 2, 5
- [61] Richard A Newcombe, Dieter Fox, and Steven M Seitz. Dynamicfusion: Reconstruction and tracking of non-rigid scenes in real-time. In *CVPR*, 2015. 1, 2
- [62] Sergio Orts-Escolano, Christoph Rhemann, Sean Fanello, Wayne Chang, Adarsh Kowdle, Yury Degtyarev, David Kim, Philip L Davidson, Sameh Khamis, Mingsong Dou, et al. Holoportation: Virtual 3d teleportation in real-time. In *UIST*, 2016. 1, 2
- [63] 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 *ICCV*, 2021. 2
- [64] 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. 2
- [65] Steven Parker, Peter Shirley, and Brian Smits. Single sample soft shadows. Technical report, Technical Report UUCS-98-019, Computer Science Department, University of Utah, 1998. 2
- [66] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 4, 5
- [67] Nikolay Patakin, Dmitry Senushkin, Anna Vorontsova, and Anton Konushin. Neural global illumination for inverse rendering. In *2023 IEEE International Conference on Image Processing (ICIP)*, pages 1580–1584. IEEE, 2023. 2
- [68] Sida Peng, Yunzhi Yan, Qing Shuai, Hujun Bao, and Xiaowei Zhou. Representing volumetric videos as dynamic mlp maps. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4252–4262, 2023. 1, 2
- [69] Eric Penner and Li Zhang. Soft 3d reconstruction for view synthesis. *ACM TOG*, 2017. 2
- [70] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2
- [71] Ruslan Rakhimov, Andrei-Timotei Ardelean, Victor Lempitsky, and Evgeny Burnaev. Npbg++: Accelerating neural point-based graphics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15969–15979, 2022. 2
- [72] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *ICCV*, pages 14335–14345, 2021. 2
- [73] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *ACM Transactions on Graphics (ToG)*, 41(4):1–14, 2022. 2
- [74] Jason Sanders and Edward Kandrot. *CUDA by example: an introduction to general-purpose GPU programming*. Addison-Wesley Professional, 2010. 5
- [75] Ruizhi Shao, Zerong Zheng, Hanzhang Tu, Boning Liu, Hongwen Zhang, and Yebin Liu. Tensor4d: Efficient neural 4d decomposition for high-fidelity dynamic reconstruction and rendering. *arXiv*, 2022. 2
- [76] Meng-Li Shih, Shih-Yang Su, Johannes Kopf, and Jia-Bin Huang. 3d photography using context-aware layered depth inpainting. In *CVPR*, 2020. 2
- [77] Dave Shreiner et al. *OpenGL programming guide: the official guide to learning OpenGL, versions 3.0 and 3.1*. Pearson Education, 2009. 4, 5
- [78] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [79] Vincent Sitzmann, Semon Rezchikov, Bill Freeman, Josh Tenenbaum, and Fredo Durand. Light field networks: Neural scene representations with single-evaluation rendering. *Advances in Neural Information Processing Systems*, 34:19313–19325, 2021. 2
- [80] Vincent Sitzmann, Justus Thies, Felix Heide, Matthias Nießner, Gordon Wetzstein, and Michael Zollhöfer. Deepvoxels: Learning persistent 3d feature embeddings. In *CVPR*, 2019.
- [81] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 2
- [82] 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. *IEEE Transactions on Visualization and Computer Graphics*, 29(5):2732–2742, 2023. 2
- [83] 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
- [84] Zhuo Su, Lan Xu, Zerong Zheng, Tao Yu, Yebin Liu, and Lu Fang. Robustfusion: Human volumetric capture with data-driven visual cues using a rgbd camera. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 246–264. Springer, 2020. 1
- [85] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8269–8279, June 2022. 2
- [86] Richard Szeliski and Polina Golland. Stereo matching with transparency and matting. In *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*, pages 517–524. IEEE, 1998. 2

- [87] Feng Wang, Sinan Tan, Xinghang Li, Zeyue Tian, and Huaping Liu. Mixed neural voxels for fast multi-view video synthesis. *arXiv preprint arXiv:2212.00190*, 2022. 2
- [88] Liao Wang, Qiang Hu, Qihan He, Ziyu Wang, Jingyi Yu, Tinne Tuytelaars, Lan Xu, and Minye Wu. Neural residual radiance fields for streamably free-viewpoint videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 76–87, 2023. 2
- [89] 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*, pages 13524–13534, 2022. 1
- [90] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P Srinivasan, Howard Zhou, Jonathan T Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4690–4699, 2021. 2, 6
- [91] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion. In *CVPR*, 2021. 2
- [92] Guanjun Wu, Taoran Yi, Jiemin Fang, Lingxi Xie, Xiaopeng Zhang, Wei Wei, Wenyu Liu, Qi Tian, and Wang Xinggang. 4d gaussian splatting for real-time dynamic scene rendering. *arXiv preprint arXiv:2310.08528*, 2023. 3
- [93] Minye Wu, Yuehao Wang, Qiang Hu, and Jingyi Yu. Multi-view neural human rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1682–1691, 2020. 2, 3, 5
- [94] Ziyi Yang, Xinyu Gao, Wen Zhou, Shaohui Jiao, Yuqing Zhang, and Xiaogang Jin. Deformable 3d gaussians for high-fidelity monocular dynamic scene reconstruction. *arXiv preprint arXiv:2309.13101*, 2023. 3
- [95] Zeyu Yang, Hongye Yang, Zijie Pan, Xiatian Zhu, and Li Zhang. Real-time photorealistic dynamic scene representation and rendering with 4d gaussian splatting. *arXiv preprint arXiv 2310.10642*, 2023. 3
- [96] Alex Yu, Sara Fridovich-Keil, Matthew Tancik, Qinlong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. *CVPR*, 2022. 2
- [97] 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*, pages 5752–5761, 2021. 1, 2, 3
- [98] Alex Yu, Vickie Ye, Matthew Tancik, and Angjoo Kanazawa. pixelNeRF: Neural radiance fields from one or few images. In *CVPR*, 2021. 2
- [99] Tao Yu, Zerong Zheng, Kaiwen Guo, Pengpeng Liu, Qionghai Dai, and Yebin Liu. Function4d: Real-time human volumetric capture from very sparse consumer rgb-d sensors. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5746–5756, 2021. 1
- [100] Tao Yu, Zerong Zheng, Kaiwen Guo, Jianhui Zhao, Qionghai Dai, Hao Li, Gerard Pons-Moll, and Yebin Liu. Doublefusion: Real-time capture of human performances with inner body shapes from a single depth sensor. In *CVPR*, 2018. 1, 2
- [101] Qiang Zhang, Seung-Hwan Baek, Szymon Rusinkiewicz, and Felix Heide. Differentiable point-based radiance fields for efficient view synthesis. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–12, 2022. 2
- [102] 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. 5
- [103] C Lawrence Zitnick, Sing Bing Kang, Matthew Uyttendaele, Simon Winder, and Richard Szeliski. High-quality video view interpolation using a layered representation. *ACM TOG*, 2004. 2