

K-Planes: Explicit Radiance Fields in Space, Time, and Appearance

Sara Fridovich-Keil*
UC Berkeley
sfk@berkeley.edu

Giacomo Meanti*
Istituto Italiano di Tecnologia
giacomo.meanti@iit.it

Frederik Rahbæk Warburg
Technical University of Denmark
frwa@dtu.dk

Benjamin Recht
UC Berkeley
brecht@berkeley.edu

Angjoo Kanazawa
UC Berkeley
kanazawa@berkeley.edu

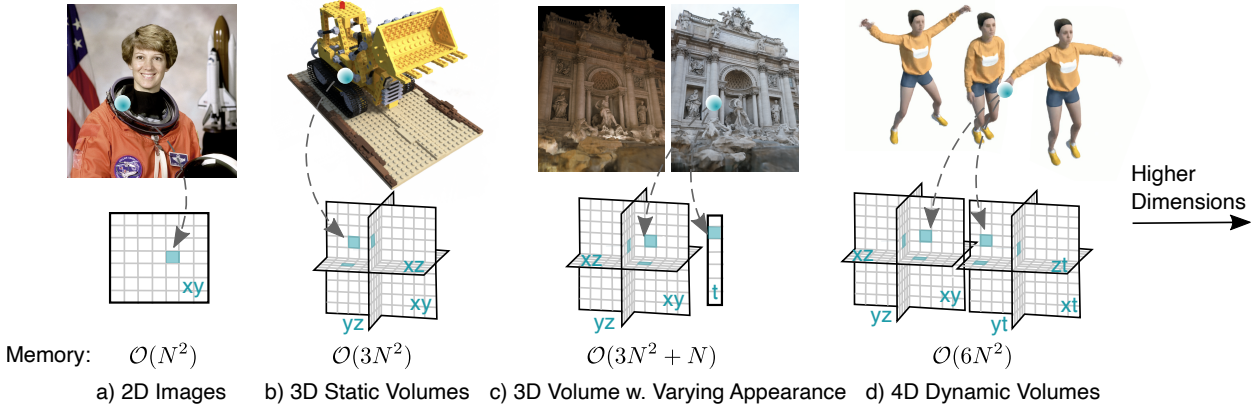


Figure 1. **Planar factorization of d -dimensional spaces.** We propose a simple planar factorization for volumetric rendering that naturally extends to arbitrary-dimensional spaces, and that scales gracefully with dimension in both optimization time and model size. We show the advantages of our approach on 3D static volumes, 3D photo collections with varying appearances, and 4D dynamic videos.

Abstract

We introduce *k*-planes, a white-box model for radiance fields in arbitrary dimensions. Our model uses $\binom{d}{2}$ (“ d -choose-2”) planes to represent a d -dimensional scene, providing a seamless way to go from static ($d = 3$) to dynamic ($d = 4$) scenes. This planar factorization makes adding dimension-specific priors easy, e.g. temporal smoothness and multi-resolution spatial structure, and induces a natural decomposition of static and dynamic components of a scene. We use a linear feature decoder with a learned color basis that yields similar performance as a nonlinear black-box MLP decoder. Across a range of synthetic and real, static and dynamic, fixed and varying appearance scenes, *k*-planes yields competitive and often state-of-the-art reconstruction fidelity with low memory usage, achieving 1000x compression over a full 4D grid, and fast optimization with a pure PyTorch implementation. For video results and code, please see sarafridov.github.io/K-Planes.

* equal contribution

1. Introduction

Recent interest in dynamic radiance fields demands representations of 4D volumes. However, storing a 4D volume directly is prohibitively expensive due to the curse of dimensionality. Several approaches have been proposed to factorize 3D volumes for static radiance fields, but these do not easily extend to higher dimensional volumes.

We propose a factorization of 4D volumes that is simple, interpretable, compact, and yields fast training and rendering. Specifically, we use six planes to represent a 4D volume, where the first three represent space and the last three represent space-time changes, as illustrated in Fig. 1(d). This decomposition of space and space-time makes our model interpretable, *i.e.* dynamic objects are clearly visible in the space-time planes, whereas static objects only appear in the space planes. This interpretability enables dimension-specific priors in time and space.

More generally, our approach yields a straightforward, prescriptive way to select a factorization of any dimension with 2D planes. For a d -dimensional space, we use $k = \binom{d}{2}$

(“ d -choose-2”) k -planes, which represent every pair of dimensions — for example, our model uses $\binom{4}{2} = 6$ *hex-planes* in 4D and reduces to $\binom{3}{2} = 3$ *tri-planes* in 3D. Choosing any other set of planes would entail either using more than k planes and thus occupying unnecessary memory, or using fewer planes and thereby forfeiting the ability to represent some potential interaction between two of the d dimensions. We call our model k -planes; Fig. 1 illustrates its natural application to both static and dynamic scenes.

Most radiance field models entail some black-box components with their use of MLPs. Instead, we seek a simple model whose functioning can be inspected and understood. We find two design choices to be fundamental in allowing k -planes to be a white-box model while maintaining reconstruction quality competitive with or better than previous black-box models [16, 30]: (1) Features from our k -planes are *multiplied* together rather than added, as was done in prior work [5, 6], and (2) our linear feature decoder uses a learned basis for view-dependent color, enabling greater adaptivity including the ability to model scenes with variable appearance. We show that an MLP decoder can be replaced with this linear feature decoder only when the planes are multiplied, suggesting that the commonly-used MLP decoder is involved in both view-dependent color and determining spatial structure.

Our factorization of 4D volumes into 2D planes leads to a high compression level without relying on MLPs, using 200 MB to represent a 4D volume whose direct representation at the same resolution would require more than 300 GB, a compression rate of three orders of magnitude. Furthermore, despite not using any custom CUDA kernels, k -planes trains orders of magnitude faster than prior implicit models and on par with concurrent hybrid models.

In summary, we present the first white-box, interpretable model capable of representing radiance fields in arbitrary dimensions, including static scenes, dynamic scenes, and scenes with variable appearance. Our k -planes model achieves competitive performance across reconstruction quality, model size, and optimization time across these varied tasks, without any custom CUDA kernels.

2. Related Work

K -planes is an interpretable, explicit model applicable to static scenes, scenes with varying appearances, and dynamic scenes, with compact model size and fast optimization time. Our model is the first to yield all of these attributes, as illustrated in Tab. 1. We further highlight that k -planes satisfies this in a simple framework that naturally extends to arbitrary dimensions.

Spatial decomposition. NeRF [24] proposed a fully implicit model with a large neural network queried many times

	Static	Appearance	Dynamic	Fast	Compact	Explicit
NeRF	✓	✗	✗	✗	✓	✗
NeRF-W	✓	✓	✗	✗	✓	✗
DVGO	✓	✗	✗	✓	✗	✗
Plenoxels	✓	✗	✗	✓	✗	✓
Instant-NGP, TensorRF	✓	✗	✗	✓	✓	✗ ²
DyNeRF, D-NeRF	–	✗	✓	✗	✓	✗
TiNeu Vox ¹ , Tensor4D ¹	–	✗	✓	✓	✓	✗
Mix Voxels ¹ , V4D ¹	–	✗	✓	✓	✗	✗
NeRFPlayer ¹	–	✗	✓	✓	✓ ³	✗
K -planes hybrid (Ours)	✓	✓	✓	✓	✓	✗
K -planes explicit (Ours)	✓	✓	✓	✓	✓	✓

¹ Recent work (within the last few months) ² TensorRF offers both hybrid and explicit versions, with a small quality gap ³ NerfPlayer offers models at different sizes, the smallest of which has < 100 million parameters but the largest of which has > 300 million parameters

Table 1. **Related work overview.** We present a simple decomposition that works for a diverse set of scenes and tasks (static, varying appearance, and dynamic). It has a low memory usage (compact) and fast training and inference time (fast). Here “fast” includes any model that can optimize within a few (< 6) hours on a single GPU, and “compact” denotes models that use less than roughly 100 million parameters. “Explicit” denotes white-box models that do not rely on MLPs.

during optimization, making it slow and essentially a black-box. Several works have used geometric representations to reduce the optimization time. Plenoxels [10] proposed a fully explicit model with trilinear interpolation in a 3D grid, which reduced the optimization time from hours to a few minutes. However, their explicit grid representation of 3D volumes, and that of DVGO [33], grows exponentially with dimension, making it challenging to scale to high resolution and completely intractable for 4D dynamic volumes.

Hybrid methods [6, 25, 33] retain some explicit geometric structure, often compressed by a spatial decomposition, alongside a small MLP feature decoder. Instant-NGP [25] proposed a multiresolution voxel grid encoded implicitly via a hash function, allowing fast optimization and rendering with a compact model. TensorRF [6] achieved similar model compression and speed by replacing the voxel grid with a tensor decomposition into planes and vectors. In a generative setting, EG3D [5] proposed a similar spatial decomposition into three planes, whose values are added together to represent a 3D volume.

Our work is inspired by the explicit modeling of Plenoxels as well as these spatial decompositions, particularly the

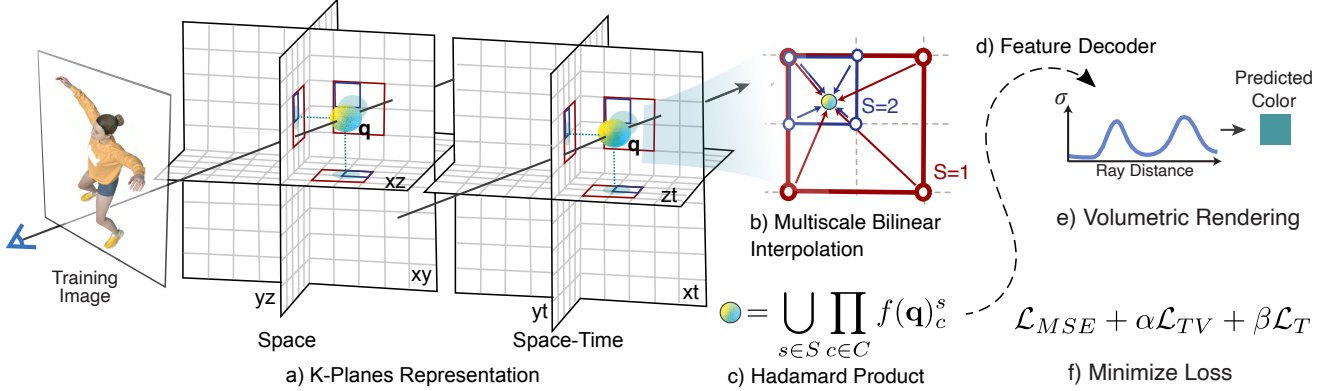


Figure 2. **Method overview.** (a) Our k -planes representation factorizes 4D dynamic volumes into six planes, three for space and three for spatiotemporal variations. To obtain the value of a 4D point $\mathbf{q} = (x, y, z, t)$, we first project the point into each plane, in which we (b) do multiscale bilinear interpolation. (c) The interpolated values are multiplied and then concatenated over the S scales. (d) These features are decoded either with a small MLP or our explicit linear decoder. (e) We follow the standard volumetric rendering formula to predict ray color and density. The model is optimized by (f) minimizing the reconstruction loss with simple regularization in space and time.

triplane model of [5], the tensor decomposition of [6], and the multiscale grid model of [25]. We also draw inspiration from Extreme MRI [26], which uses a multiscale low-rank decomposition to represent 4D dynamic volumes in magnetic resonance imaging. These spatial decomposition methods have been shown to offer a favorable balance of memory efficiency and optimization time for static scenes. However, it is not obvious how to extend these factorizations to 4D volumes in a memory-efficient way. K -planes defines a unified framework that enables efficient and interpretable factorizations of 3D and 4D volumes and trivially extends to even higher dimensional volumes.

Dynamic volumes. Applications such as Virtual Reality (VR) and Computed Tomography (CT) often require the ability to reconstruct 4D volumes. Several works have proposed extensions of NeRF to dynamic scenes. The two most common schemes are (1) modeling a deformation field on top of a static *canonical* field [8, 9, 17, 27, 30, 36, 43], or (2) directly learning a radiance field conditioned on time [12, 16, 17, 28, 41]. The former makes decomposing static and dynamic components easy [40, 43], but struggles with changes in scene topology (e.g. when a new object appears), while the latter makes disentangling static and dynamic objects hard. A third strategy is to choose a compact 3D spatial representation and essentially repeat it at each timestep (e.g. NeRFPlayer [32]), resulting in a model that ignores interactions between space and time and can become impractically large for long videos.

Further, some of these models are fully implicit [16, 30] and thus suffer from extremely long training times (e.g. DyNeRF used 8 GPUs for 1 week to train a single scene), as well as being completely black-box. Others use partially explicit decompositions for video [9, 11, 14, 18, 19, 31, 32, 37], usually combining some voxel or spatially decomposed fea-

ture grid with one or more MLP components for feature decoding and/or representing scene dynamics. Most closely related to k -planes is Tensor4D [31], which uses 9 planes to decompose 4D volumes. K -planes is less redundant (e.g. Tensor4D includes two yt planes), does not rely on multiple MLPs, and offers a simpler factorization that naturally generalizes to static and dynamic scenes. Our method combines a fully explicit representation with a built-in decomposition of static and dynamic components, the ability to handle arbitrary topology and lighting changes over time, fast optimization, and compactness.

Appearance embedding. Reconstructing large environments from photographs taken with varying illumination is another domain in which implicit methods have shown appealing results, but hybrid and explicit approaches have not yet gained a foothold. NeRF-W [20] was the first to demonstrate photorealistic view synthesis in such environments. They augment a NeRF-based model with a learned global appearance code per frame, enabling it to explain away changes in appearance, such as time of day. With Generative Latent Optimization (GLO) [4], these appearance codes can further be used to manipulate the scene appearance by interpolation in the latent appearance space. Block-NeRF [34] employs similar appearance codes.

We show that our k -planes representation can also effectively reconstruct these unbounded environments with varying appearance. We similarly extend our model – either the learned color basis in the fully explicit version, or the MLP decoder in the hybrid version – with a global appearance code to disentangle global appearance from a scene without affecting geometry. To the best of our knowledge, ours is both the first fully explicit and the first hybrid method to successfully reconstruct these challenging scenes.

3. K-planes model

We propose a simple and interpretable model for representing scenes in arbitrary dimensions. Our representation yields low memory usage and fast training and rendering. The k -planes factorization, illustrated in Fig. 2, models a d -dimensional scene using $k = \binom{d}{2}$ planes representing every combination of two dimensions. For example, for static 3D scenes, this results in *tri-planes* with $\binom{3}{2} = 3$ planes representing xy , xz , and yz . For dynamic 4D scenes, this results in *hex-planes*, with $\binom{4}{2} = 6$ planes including the three space-only planes and three space-time planes xt , yt , and zt . Should we wish to represent a 5D space, we could use $\binom{5}{2} = 10$ *deca-planes*. In the following section, we describe the 4D instantiation of our k -planes factorization.

3.1. Hex-planes

The Hex-planes factorization uses six planes. We refer to the space-only planes as \mathbf{P}_{xy} , \mathbf{P}_{xz} , and \mathbf{P}_{yz} , and the space-time planes as \mathbf{P}_{xt} , \mathbf{P}_{yt} , and \mathbf{P}_{zt} . Assuming symmetric spatial and temporal resolution N for simplicity of illustration, each of these planes has shape $N \times N \times M$, where M is the size of stored features that capture the density and view-dependent color of the scene.

We obtain the features of a 4D coordinate $\mathbf{q} = (i, j, k, \tau)$ by normalizing it between $[0, N)$ and projecting it onto these six planes

$$f(\mathbf{q})_c = \psi(\mathbf{P}_c, \pi_c(\mathbf{q})), \quad (1)$$

where π_c projects \mathbf{q} onto the c 'th plane and ψ denotes bilinear interpolation of a point into a regularly spaced 2D grid. We repeat Eq. (1) for each plane $c \in C$ to obtain feature vectors $f(\mathbf{q})_c$. We combine these features over the six planes using the Hadamard product (elementwise multiplication) to produce a final feature vector of length M

$$f(\mathbf{q}) = \prod_{c \in C} f(\mathbf{q})_c. \quad (2)$$

These features will eventually be decoded into color and density using either a linear decoder or an MLP, described in Sec. 3.3.

Why Hadamard product? In 3D, k -planes reduces to the tri-plane factorization, which is similar to [5] except that the elements are multiplied. A natural question is why we multiply rather than add, as has been used in prior work with tri-plane models [5, 29]. Fig. 3 illustrates that combining the planes by multiplication allows k -planes to produce spatially localized signals, which is not possible with addition.

This selection ability of the Hadamard product produces substantial rendering improvements for linear decoders and modest improvement for MLP decoders, as shown in Tab. 2.

This suggests that the MLP decoder is involved in both view-dependent color and determining spatial structure. The Hadamard product relieves the feature decoder of this extra task and makes it possible to reach similar performance using a linear decoder solely responsible for view-dependent color.

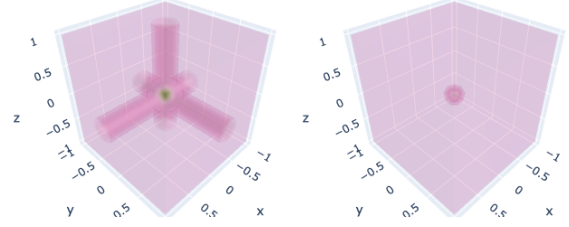


Figure 3. **Addition versus Hadamard product.** Elementwise addition of plane features (left) compared to multiplication (right), illustrated in a 3D triplane example. In both cases a single entry in each plane is positive and the rest are zero, selecting a single 3D point by multiplication but producing intersecting lines by addition. This selection ability of multiplication improves the expressivity of our explicit model.

Plane Combination	Explicit	Hybrid	# params ↓
Multiplication	36.36	36.97	34M
Addition	30.59	36.73	34M

Table 2. **Ablation study over Hadamard product.** Multiplication of plane features yields a large improvement in PSNR ↑ for our explicit model, whereas our hybrid model can use its MLP decoder to compensate for the less expressive addition of planes. This experiment uses the static *Lego* scene [24] with 4 scales: 64, 128, 256, and 512, and 32 features per scale.

3.2. Interpretability

The separation of space-only and space-time planes makes the model interpretable and enables us to incorporate dimension-specific priors. For example, if a region of the scene never moves, its temporal component will always be 1 (the multiplicative identity), thereby just using the features from the space planes. This offers compression benefits since a static region can easily be identified and compactly represented. Furthermore, the space-time separation improves interpretability, *i.e.* we can track the changes in time by visualizing the elements in the time-space planes that are not 1. This simplicity, separation, and interpretability make adding priors straightforward.

Multiscale planes. To encourage spatial smoothness and coherence, our model contains multiple copies at differ-

ent spatial resolutions, for example 64, 128, 256, and 512. Models at each scale are treated separately, and the M -dimensional feature vectors from different scales are concatenated together before being passed to the decoder. The red and blue squares in Fig. 2a-b illustrate bilinear interpolation with multiscale planes. Inspired by the multiscale hash mapping of Instant-NGP [25], this representation efficiently encodes spatial features at different scales, allowing us to reduce the number of features stored at the highest resolution and thereby further compressing our model. Empirically, we do not find it necessary to represent our time dimension at multiple scales.

Total variation in space. Spatial total variation regularization encourages sparse gradients, encoding the prior that edges are sparse in space. We encourage this in 1D over the spatial dimensions of each of our space-time planes and in 2D over our space-only planes:

$$\mathcal{L}_{TV}(\mathbf{P}) = \frac{1}{|C|n^2} \sum_{c,i,j} (\|\mathbf{P}_c^{i,j} - \mathbf{P}_c^{i-1,j}\|_2^2 + \|\mathbf{P}_c^{i,j} - \mathbf{P}_c^{i,j-1}\|_2^2), \quad (3)$$

where i, j are indices on the plane’s resolution. Total variation is a common regularizer in inverse problems and was used in Plenoxels [10] and TensorRF [6].

Smoothness in time. We encourage smooth motion with a 1D Laplacian (second derivative) filter

$$\mathcal{L}_{smooth}(\mathbf{P}) = \frac{1}{|C|n^2} \sum_{c,i,t} \|\mathbf{P}_c^{i,t-1} - 2\mathbf{P}_c^{i,t} + \mathbf{P}_c^{i,t+1}\|_2^2, \quad (4)$$

to penalize sharp “acceleration” over time. We only apply this regularizer on the time dimension of our space-time planes. Please see the appendix for an ablation study.

Sparse transients. We want the static part of the scene to be modeled by the space-only planes. We encourage this separation of space and time by initializing the features in the space-time planes as 1 (the multiplicative identity) and using an ℓ_1 regularizer on these planes during training:

$$\mathcal{L}_{sep}(\mathbf{P}) = \sum_c \|\mathbf{1} - \mathbf{P}_c\|_1, \quad c \in \{xt, yt, zt\}. \quad (5)$$

In this way, the space-time plane features of the k -planes decomposition will remain fixed at 1 if the corresponding spatial content does not change over time.

3.3. Feature decoders

We offer two methods to decode the M -dimensional temporally- and spatially-localized feature vector $f(\mathbf{q})$ from Eq. (2) into density, σ , and view-dependent color, \mathbf{c} .

Learned color basis: a linear decoder and explicit model.

Plenoxels [10], Plenotrees [42], and TensorRF [6] proposed model versions where spatially-localized features are interpreted as coefficients of the spherical harmonic basis functions and used to describe view-dependent color. Such spherical harmonic decoders can offer both high-fidelity reconstructions and enhanced interpretability compared to MLP decoders. However, spherical harmonic coefficients are difficult to optimize, and their expressivity is limited by the number of spherical harmonic basis functions used. Often only harmonics up to degree two are used, thereby producing low-frequency (blurry) specular reflections.

Instead, we replace the spherical harmonic basis functions with a learned basis, retaining the interpretability of treating features as coefficients for a linear decoder yet increasing the expressivity of the basis and allowing it to adapt to each scene, as was proposed in NeX [39]. We represent the basis as a small MLP that predicts a red $b_R(\mathbf{d}) \in \mathbb{R}^M$, green $b_G(\mathbf{d}) \in \mathbb{R}^M$, and blue $b_B(\mathbf{d}) \in \mathbb{R}^M$ basis whose input is view direction \mathbf{d} . The MLP serves as an adaptive drop-in replacement for the spherical harmonic basis functions repeated over the three color channels. We obtain the color values

$$RGB(\mathbf{q}, \mathbf{d}) = \bigcup_{i \in \{R, G, B\}} f(\mathbf{q}) \cdot b_i(\mathbf{q}, \mathbf{d}), \quad (6)$$

where \cdot denotes the dot product and \cup denotes concatenation (union). Similarly, we optimize a linear decoder for density by predicting a $b_\sigma \in \mathbb{R}^M$ density basis, independent of the view direction:

$$\sigma(\mathbf{q}) = f(\mathbf{q}) \cdot b_\sigma(\mathbf{q}), \quad (7)$$

where $\sigma(\mathbf{q})$ is the density of a d -dimensional point. The predicted color and density values are finally passed through a nonlinear function to ensure they are in the valid ranges: we apply sigmoid to the color values to ensure they lie in $[0, 1]$, and exponential (with truncated gradient) to the density values to ensure they are nonnegative.

MLP decoder: a hybrid model.

Our model can also be used with an MLP decoder like that of Instant-NGP [25] and DVGO [33], turning it into a hybrid model. In this version, features are decoded by two small MLPs, one g_σ that maps the spatially-localized features into density σ and additional features \hat{f} , and another g_{RGB} that maps these additional features as well as the embedded view direction $\gamma(\mathbf{d})$ into RGB color

$$\begin{aligned} \sigma(\mathbf{q}), \hat{f}(\mathbf{q}) &= g_\sigma(f(\mathbf{q})) \\ RGB(\mathbf{q}, \mathbf{d}) &= g_{RGB}(\hat{f}(\mathbf{q}), \gamma(\mathbf{d})). \end{aligned} \quad (8)$$

As in the linear decoder case, the predicted density and color values are finally normalized via exponential (with truncated gradient) and sigmoid, respectively.

Global appearance. We also show a simple extension of our k -planes model that enables it to represent scenes with consistent, static geometry viewed under varying lighting or appearance conditions. Such scenes appear in the Phototourism [15] dataset of famous landmarks photographed by tourists at different times of day and in different weather. To model this variable appearance, we augment k -planes with a matrix of shape $M \times T$, an M -vector of features for each training image index $1 \dots T$. Similar to NeRF-W [20], we optimize this per-image feature vector and pass it as an additional input to either our MLP learned color basis b_R, b_G, b_B , in our explicit version, or to our MLP color decoder g_{RGB} , in our hybrid version, so that in either case it can affect color but not geometry.

3.4. Optimization details

Full details of our optimization may be found in our released code. Here we specify three components we modify from prior work to improve flexibility and speed.

Contraction and normalized device coordinates. For forward-facing scenes, we apply normalized device coordinates (NDC) [24] to better allocate our resolution while enabling unbounded depth. We also implement an ℓ_∞ version (rather than ℓ_2) of the scene contraction proposed in Mip-NeRF 360 [2], which we use on the unbounded Phototourism scenes.

Proposal sampling. We use a variant of the proposal sampling strategy from Mip-NeRF 360 [2], with a small instance of k -planes as each density model. Proposal sampling works by first using a fixed sampling strategy along each ray into a density model, and then adaptively choosing samples based on the densities of these initial samples. We use a two-stage proposal sampler, resulting in both fewer samples that must be evaluated in our full model and in sharper details by placing those samples closer to object surfaces. The density models used for proposal sampling are trained with the histogram loss [2].

Importance sampling. For multiview dynamic scenes, we implement a version of the importance sampling based on temporal difference (IST) strategy from DyNeRF [16]. During the last portion of optimization, we sample training rays according to the maximum absolute variation in their color within 25 frames, slightly less than one second before or after. This results in higher sampling probabilities in the dynamic region. We apply this strategy after

the static scene has converged through standard optimization with uniformly sampled rays. In our experiments, this sampling strategy has only a modest impact on full-frame metrics but improves visual quality in the small dynamic region. Note that importance sampling cannot be used for monocular videos or any datasets with moving cameras.

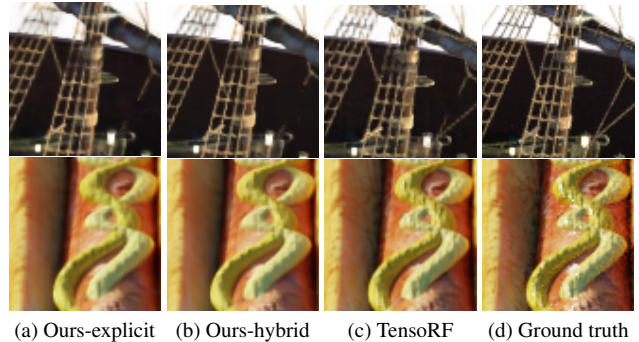


Figure 4. **Zoomed qualitative results on static NeRF scenes.** Visual comparison of k -planes, TensorRF [6], and the ground truth, on *ship* (top) and *hotdog* (bottom).

4. Results

We demonstrate the broad applicability of our planar decomposition via experiments in three domains: static scenes (both bounded 360° and unbounded forward-facing), dynamic scenes (forward-facing multi-view and bounded 360° monocular), and Phototourism scenes with variable appearance. For all experiments, we report the metrics PSNR (pixel-level similarity) and SSIM¹ [38] (structural similarity), as well as approximate training time and number of parameters (in millions), in Tab. 3. Blank entries in Tab. 3 denote baseline methods for which the corresponding information is not readily available. Full per-scene results may be found in the appendix.

4.1. Static scenes

We first demonstrate our triplane model on the bounded, 360° , synthetic scenes from NeRF [24]. We use a model with four symmetric spatial resolutions $N \in \{64, 128, 256, 512\}$ and feature length $M = 32$ at each scale; please see the appendix for ablation studies over these hyperparameters. The explicit version of our model matches the prior state-of-the-art in terms of quality metrics, while the hybrid version achieves slightly higher quality metrics. Fig. 4 shows zoomed-in visual results on a small sampling of scenes.

¹Note that among prior work, some evaluate using a public implementation of SSIM from MipNeRF [1] whereas others use the implementation in scikit-image, which tends to produce higher values. For fair comparison on each dataset we make a best effort to use the same SSIM implementation as the relevant prior work.

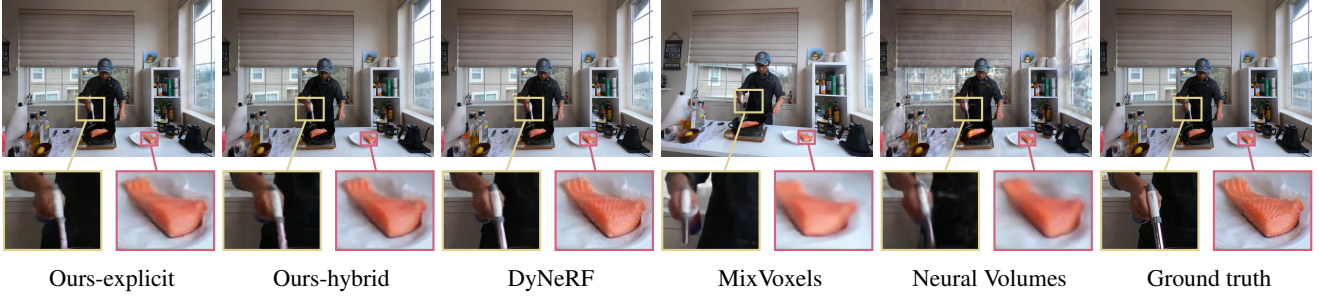


Figure 5. **Qualitative video results.** Our hexplane model rivals the rendering quality of state-of-the-art neural rendering methods. Our renderings were obtained after at most 4 hours of optimization on a single GPU whereas DyNeRF trained for a week on 8 GPUs. MixVoxels frame comes from a slightly different video rendering, and is thus slightly shifted.

We also present results of our triplane model on the unbounded, forward-facing, real scenes from LLFF [23]. Our results on this dataset are similar to the synthetic static scenes; both versions of our model match or exceed the prior state-of-the-art, with the hybrid version achieving slightly higher metrics than the fully explicit version. Fig. 6 shows zoomed-in visual results on a small sampling of scenes.

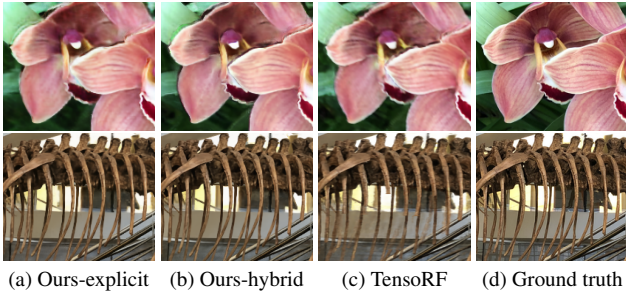


Figure 6. **Zoomed qualitative results on static LLFF scenes** [23]. Visual comparison of k -planes, TensoRF [6], and the ground truth, on *orchids* (top) and *T-rex* (bottom).

4.2. Dynamic scenes

We evaluate our hexplane model on two dynamic scene datasets: a set of synthetic, bounded, 360°, monocular videos from D-NeRF [30] and a set of real, unbounded, forward-facing, multiview videos from DyNeRF [16].

The D-NeRF dataset contains eight videos of varying duration, from 50 frames to 200 frames per video. Each timestep has a single training image from its own camera viewpoint; the camera “teleports” between adjacent timesteps [13]. Standardized test views are from novel camera positions at a range of timestamps throughout the video. Both our explicit and hybrid models outperform D-NeRF in both quality metrics and training time (by a large margin), though they do not surpass very recent hybrid methods TiNeuVox [9] and V4D [11], as shown in Fig. 7.

The DyNeRF dataset contains six 10-second videos recorded at 30 fps simultaneously on approximately 15-20

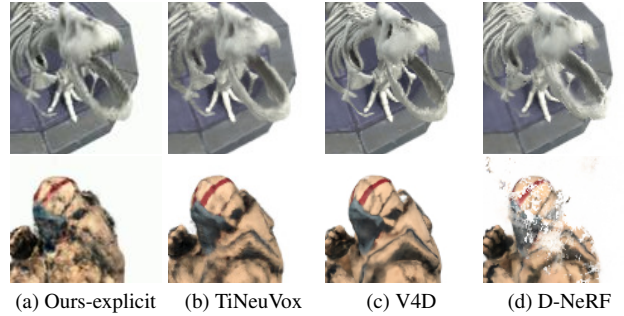


Figure 7. **Zoomed qualitative results on monocular dynamic scenes from D-NeRF** [30]. Visual comparison of k -planes, D-NeRF [30], TiNeuVox [9] and V4D [11], on *t-rex* (top) and *hook* (bottom).

video cameras from a range of forward-facing view directions; the exact number of cameras varies per scene because a few cameras produced miscalibrated videos. A central camera is reserved for evaluation, and training uses frames from the remaining cameras. Both our methods again produce similar quality metrics to prior state-of-the-art, including a very recent hybrid method MixVoxels [37], with our hybrid method achieving higher quality metrics. Please see Fig. 5 for a visual comparison.

4.2.1 Decomposing time and space

One neat consequence of our planar decomposition of time and space is that it naturally disentangles dynamic and static portions of the scene. The static-only part of the scene can be obtained by setting the three time planes to one (the multiplicative identity). Subtracting the static-only rendered image from the full rendering (i.e. with the time plane parameters not set to 1), we can reveal the dynamic part of the scene. Fig. 8 shows this decomposition of time and space. This natural volumetric disentanglement of a scene into static and dynamic regions may enable many applications across augmented and virtual reality [3].

We can also visualize the time planes to better under-

	PSNR \uparrow	SSIM \uparrow	Train Time \downarrow	# Params \downarrow
NeRF [24] (static, synthetic)				
Ours-explicit	33.13	0.964	38 min	34M
Ours-hybrid	33.62	0.967	38 min	34M
Plenoxels [10]	31.71	0.958	11 min	~500M
TensorRF [6]	33.14	0.963	17 min	18M
I-NGP [25]	33.18	-	5 min	~16M
LLFF [23] (static, real)				
Ours-explicit	26.78	0.841	33 min	19M
Ours-hybrid	26.92	0.847	33 min	19M
Plenoxels	26.29	0.839	24 min	~500M
TensorRF	26.73	0.839	25 min	45M
D-NeRF [30] (dynamic, synthetic)				
Ours-explicit	30.39	0.96	52 min	37M
Ours-hybrid	30.84	0.96	52 min	37M
D-NeRF	29.67	0.95	48 hrs	1-3M
TiNeuVox ¹ [9]	32.67	0.97	30 min	~12M
V4D ¹ [11]	33.72	0.98	4.9 hrs	275M
DyNeRF [16] (dynamic, real)				
Ours-explicit	30.88	0.960	3.7 hrs	51M
Ours-hybrid	31.63	0.964	1.8 hrs	27M
DyNeRF [16]	² 29.58	-	1344 hrs	7M
LLFF [23]	² 23.24	-	-	-
MixVoxels-L ¹ [37]	30.80	0.960	1.3 hrs	125M
Phototourism [15] (variable appearance)				
Ours-explicit	22.25	0.859	35 min	36M
Ours-hybrid	22.92	0.877	35 min	36M
NeRF-W [20]	27.00	0.962	384 hrs	~2M
NeRF-W (public) ³	19.70	0.764	164 hrs	~2M
LearnIt [35]	19.26	-	-	-

¹ MixVoxels, TiNeuVox and V4D were released within the last few months. TiNeuVox uses half-resolution images for training and evaluation (V4D code is not yet public, so their resolution is unknown). ² DyNeRF and LLFF only report metrics on the *flame salmon* video; average performance may be higher as this is one of the more challenging videos. ³ Open-source version https://github.com/kweal23/nerf_pl/tree/nerfw where we re-implement test-time optimization as for *k*-planes.

Table 3. **Results.** *First section:* average over the 8 synthetic static scenes from NeRF [24]; *Second section:* average over the 8 real forward-facing static scenes from LLFF [23]; *Third section:* average over the 8 synthetic monocular “teleporting camera” dynamic scenes from D-NeRF [30]; *Fourth section:* average over the 6 real multiview forward-facing dynamic scenes from DyNeRF [16]; *Fifth section:* average over three scenes from the Phototourism dataset [15]; note these scenes use MS-SSIM (multiscale structural similarity) instead of SSIM. *K*-planes timing results are based on a single NVIDIA A30 GPU. Please see the appendix for per-scene results and the website for video reconstructions.

stand where motion occurs in a video. Fig. 9 shows the averaged features learned by the *xt* plane in our model for the *flame salmon* and *cut beef* DyNeRF videos, in which

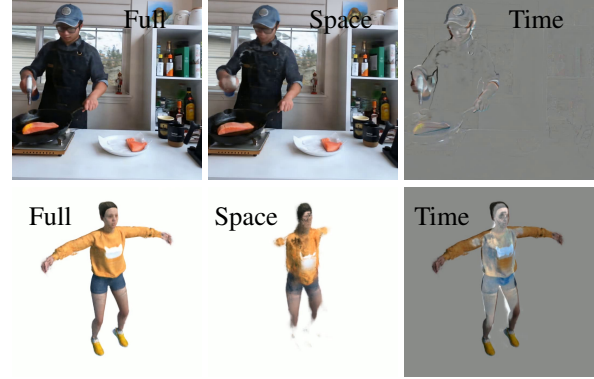


Figure 8. **Decomposition of space and time.** *K*-planes (left) naturally decomposes a 3D video into static and dynamic components. We render the static part (middle) by setting the time planes to the identity, and the remainder (right) is the dynamic part. Top shows the *flame salmon* multiview video [16] and bottom shows the *jumping jacks* monocular video [30].

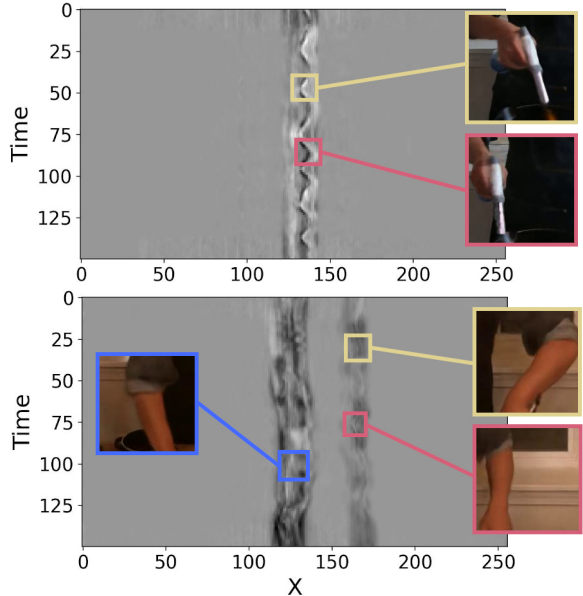


Figure 9. **Visualization of a time plane.** The *xt* plane highlights the dynamic regions in the scene. The wiggly patterns across time correspond to the motion of the person’s hands and cooking tools, in the *flame salmon* scene (top) where only one hand moves and the *cut beef* scene (bottom) where both hands move.

we can identify the motions of the hands in both space and time. The *xt* plane learns to be sparse, with most entries equal to the multiplicative identity, due to a combination of our sparse transients prior and the true sparsity of motion in the video. For example, in the upper portion of Fig. 8 one of the cook’s arms contains most of the motion, while in the lower figure both arms are moving. Having access to such an explicit representation of time allows us to add time-specific priors.

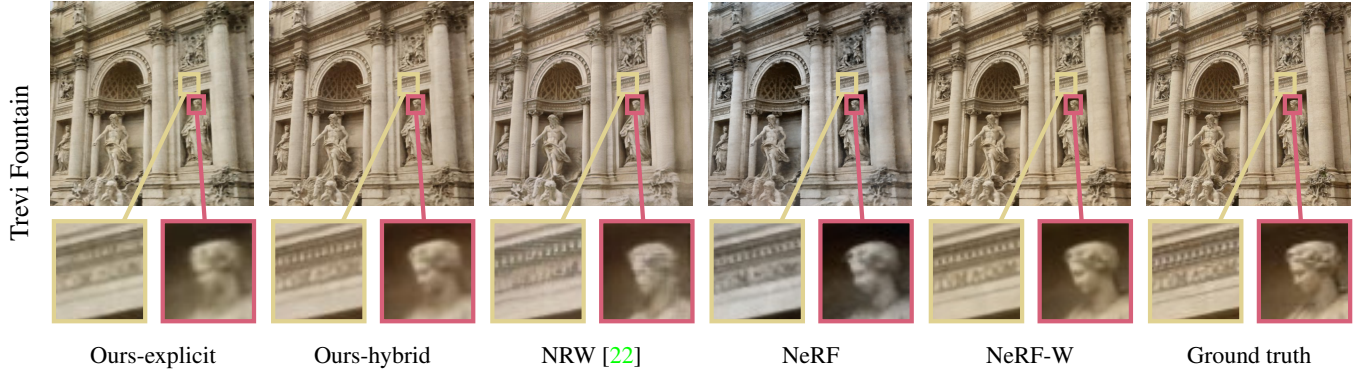


Figure 10. **Qualitative results from Phototourism dataset.** We compare our model with strong baselines. Our method captures the geometry and appearance of the scene, but produces slightly lower resolution results than NeRF-W. Note that our model optimizes in just 35 minutes on a single GPU compared to NeRF-W, which takes 2 days on 8 GPUs.

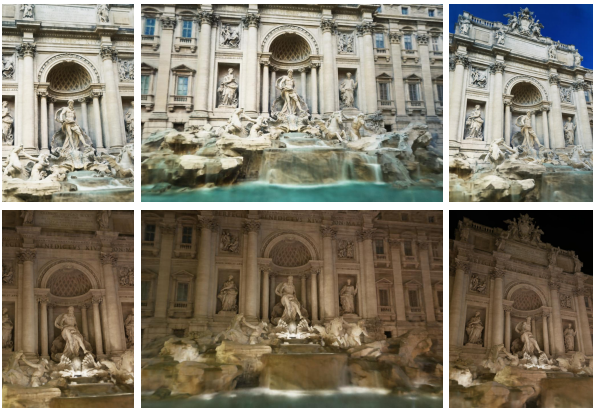


Figure 11. **Appearance interpolation.** Like NeRF-W [20], we can interpolate our appearance code to alter the visual appearance of landmarks. We show three test views from the *Trevi fountain* with appearance codes corresponding to day and night.

4.3. Variable appearance

Our variable appearance experiments use the Phototourism dataset [15], which includes photos of well-known landmarks taken by tourists with arbitrary view directions, lighting conditions, and transient occluders, mostly other tourists. Our experimental conditions parallel those of NeRF-W [20]: we train on more than a thousand tourist photographs and test on a standard set that is free of transient occluders. Like NeRF-W, we evaluate on test images by optimizing our per-image appearance feature on the left half of the image and computing metrics on the right half. Results are visualized in Fig. 10.

Also similar to NeRF-W [4, 20], we can interpolate in the appearance code space. Since only the color decoder (and not the density decoder) takes the appearance code as input, our approach is guaranteed not to change the geometry, regardless of whether we use our explicit or our hybrid model. Fig. 11 shows that our planar decomposition with a

32-dimensional appearance code is sufficient to accurately capture global appearance changes in the scene.

5. Conclusions

We introduced a simple yet versatile method to decompose a d -dimensional space into $\binom{d}{2}$ planes, which can be optimized directly from indirect measurements and scales gracefully in model size and optimization time with increasing dimension, without any custom CUDA kernels. We demonstrated that the proposed k -planes decomposition applies naturally to reconstruction of static 3D scenes as well as dynamic 4D videos, and with the addition of a global appearance code can also extend to the more challenging task of unconstrained scene reconstruction. K -planes is the first explicit, simple model to demonstrate competitive performance across such varied tasks.

Acknowledgments. Many thanks to Matthew Tancik, Ruilong Li, and other members of KAIR for helpful discussion and pointers to resources. We also thank the DyNeRF authors for their response to our questions about their method.

References

- [1] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 5835–5844. IEEE, 2021. 6
- [2] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 5460–5469. IEEE, 2022. 6
- [3] Sagie Benaim, Frederik Warburg, Peter Ebert Christensen, and Serge Belongie. Volumetric disentanglement for 3d scene manipulation, 2022. 7
- [4] Piotr Bojanowski, Armand Joulin, David Lopez-Paz, and Arthur Szlam. Optimizing the latent space of generative networks, 2017. 3, 9
- [5] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas J. Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3d generative adversarial networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 16102–16112. IEEE, 2022. 2, 3, 4
- [6] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision (ECCV)*, 2022. 2, 3, 5, 6, 7, 8, 13
- [7] Boyang Deng, Jonathan T. Barron, and Pratul P. Srinivasan. JaxNeRF: an efficient JAX implementation of NeRF, 2020. 13
- [8] Yilun Du, Yinan Zhang, Hong-Xing Yu, Joshua B. Tenenbaum, and Jiajun Wu. Neural radiance flow for 4d view synthesis and video processing. In *ICCV*, pages 14304–14314, 2021. 3
- [9] Jiemin Fang, Taoran Yi, Xinggang Wang, Lingxi Xie, Xiaopeng Zhang, Wenyu Liu, Matthias Nießner, and Qi Tian. Fast dynamic radiance fields with time-aware neural voxels. In *SIGGRAPH Asia 2022 Conference Papers*. ACM, nov 2022. 3, 7, 8, 14
- [10] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 5491–5500. IEEE, 2022. 2, 5, 8, 13
- [11] Wanshui Gan, Hongbin Xu, Yi Huang, Shifeng Chen, and Naoto Yokoya. V4d: Voxel for 4d novel view synthesis, 2022. 3, 7, 8, 14
- [12] Chen Gao, Ayush Saraf, Johannes Kopf, and Jia-Bin Huang. Dynamic view synthesis from dynamic monocular video. In *ICCV*, pages 5712–5721, 2021. 3
- [13] Hang Gao, Ruilong Li, Shubham Tulsiani, Bryan Russell, and Angjoo Kanazawa. Monocular dynamic view synthesis: A reality check. In *NeurIPS*, 2022. 7
- [14] Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. Neural deformable voxel grid for fast optimization of dynamic view synthesis. In *Proceedings of the Asian Conference on Computer Vision (ACCV)*, 2022. 3
- [15] Yuhe Jin, Dmytro Mishkin, Anastasiia Mishchuk, Jiri Matas, Pascal Fua, Kwang Moo Yi, and Eduard Trulls. Image matching across wide baselines: From paper to practice. *Int. J. Comput. Vis.*, 129(2):517–547, 2021. 6, 8, 9
- [16] Tianye Li, Mira Slavcheva, Michael Zollhoefer, Simon Green, Christoph Lassner, Changil Kim, Tanner Schmidt, Steven Lovegrove, Michael Goesele, Richard Newcombe, and Zhaoyang Lv. Neural 3d video synthesis from multi-view video. In *CVPR*, pages 5511–5521, 2022. 2, 3, 6, 7, 8, 14
- [17] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. In *CVPR*, 2021. 3
- [18] Jia-Wei Liu, Yan-Pei Cao, Weijia Mao, Wenqiao Zhang, David Junhao Zhang, Jussi Keppo, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Devrf: Fast deformable voxel radiance fields for dynamic scenes. *arxiv:2205.15723*, 2022. 3
- [19] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019. 3
- [20] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. NeRF in the Wild: Neural Radiance Fields for Unconstrained Photo Collections. In *CVPR*, 2021. 3, 6, 8, 9, 15
- [21] N. Max. Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics*, 1(2):99–108, 1995. 12
- [22] Moustafa Meshry, Dan B. Goldman, Sameh Khamis, Hugues Hoppe, Rohit Pandey, Noah Snavely, and Ricardo Martin-Brualla. Neural rerendering in the wild. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 6878–6887. Computer Vision Foundation / IEEE, 2019. 9
- [23] 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 Trans. Graph.*, 38(4):29:1–29:14, 2019. 7, 8, 13, 14
- [24] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*, pages 405–421. Springer, 2020. 2, 4, 6, 8, 12, 13
- [25] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, July 2022. 2, 3, 5, 8, 13

- [26] Frank Ong, Xucheng Zhu, Joseph Y. Cheng, Kevin M. Johnson, Peder E. Z. Larson, Shreyas S. Vasanawala, and Michael Lustig. Extreme mri: Large-scale volumetric dynamic imaging from continuous non-gated acquisitions. *Magnetic Resonance in Medicine*, 84(4):1763–1780, 2020. 3
- [27] 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. 3
- [28] 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. *ACM Trans. Graph.*, 40(6), dec 2021. 3
- [29] Songyou Peng, Michael Niemeyer, Lars M. Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part III*, volume 12348 of *Lecture Notes in Computer Science*, pages 523–540. Springer, 2020. 4
- [30] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *CVPR*, 2021. 2, 3, 7, 8, 12, 14
- [31] 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, 2022. 3, 14
- [32] Liangchen Song, Anpei Chen, Zhong Li, Zhang Chen, Lele Chen, Junsong Yuan, Yi Xu, and Andreas Geiger. Nerf-player: A streamable dynamic scene representation with decomposed neural radiance fields, 2022. 3
- [33] Cheng Sun, Min Sun, and Hwann-Tzong Chen. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *CVPR*, 2022. 2, 5, 13
- [34] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P. Srinivasan, Jonathan T. Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis, 2022. 3
- [35] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P. Srinivasan, Jonathan T. Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 2846–2855. Computer Vision Foundation / IEEE, 2021. 8, 15
- [36] 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 *ICCV*. IEEE, 2021. 3
- [37] Feng Wang, Sinan Tan, Xinghang Li, Zeyue Tian, and Huaping Liu. Mixed neural voxels for fast multi-view video synthesis, 2022. 3, 7, 8, 14
- [38] Zhou Wang, A.C. Bovik, H.R. Sheikh, and E.P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 6
- [39] Suttisak Wizadwongsa, Pakkapon Phongthawee, Jiraphon Yenphraphai, and Supasorn Suwajanakorn. Nex: Real-time view synthesis with neural basis expansion, 2021. 5
- [40] Tianhao Wu, Fangcheng Zhong, Andrea Tagliasacchi, Forrester Cole, and Cengiz Oztireli. D²nerf: Self-supervised decoupling of dynamic and static objects from a monocular video, 2022. 3
- [41] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9421–9431, 2021. 3
- [42] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 5732–5741. IEEE, 2021. 5
- [43] Wentao Yuan, Zhaoyang Lv, Tanner Schmidt, and Steven Lovegrove. STaR: Self-supervised tracking and reconstruction of rigid objects in motion with neural rendering. In *CVPR*, 2021. 3

6. Appendix

6.1. Volumetric rendering

We use the same volume rendering formula as NeRF [24], originally from [21], where the color of a pixel is represented as a sum over samples taken along the corresponding ray through the volume:

$$\sum_{i=1}^N \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right) (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i \quad (9)$$

where the first exp represents ray transmission to sample i , $1 - \exp(-\sigma_i \delta_i)$ is the absorption by sample i , σ_i is the (post-activation) density of sample i , and \mathbf{c}_i is the color of sample i , with distance δ_i to the next sample.

6.2. Per-scene results

In this section, we provide quantitative metrics for each of the three tasks we study, for each scene individually. Tab. 7 reports metrics on the static synthetic scenes, Tab. 8 reports metrics on the static real forward-facing scenes, Tab. 9 reports metrics on the dynamic synthetic monocular “teleporting camera” scenes, Tab. 10 reports metrics on the dynamic real forward-facing multiview scenes, and Tab. 11 reports metrics on the Phototourism scenes.

6.3. Ablation studies

Multiscale. In Tab. 4, we ablate our model on the static *Lego* scene [24] with respect to our multiscale planes, to assess the value of including copies of our model at different scales.

Scales	Explicit PSNR \uparrow	Hybrid PSNR \uparrow	# params \downarrow
64, 128, 256, 512	36.36	36.97	34M
128, 256, 512	36.44	37.23	34M
256, 512	36.40	37.19	32M
512	35.65	36.92	26M
64, 128, 256	33.79	35.22	9M

Table 4. **Ablation study over scales.** Including even a single lower scale improves performance, particularly for our explicit model, with a relatively small overhead in train time and model size. Using lower scales only (excluding resolution 512³) reduces both quality (PSNR) and model size substantially. This experiment uses the static *Lego* scene; each scale uses 32 features.

Feature length. In Tab. 5, we ablate our model on the static *Lego* scene with respect to the feature dimension M learned at each scale.

Feature Length (M)	Explicit PSNR \uparrow	Hybrid PSNR \uparrow	# params \downarrow
4	32.82	34.88	5M
8	34.65	35.96	9M
16	35.57	36.78	17M
32	36.36	36.97	34M
64	36.93	37.66	67M

Table 5. **Ablation study over feature length M .** Increasing the feature length M learned at each scale consistently improves quality for both our models, with a corresponding linear increase in model size and optimization time. Our experiments in the main text use a mixture of $M = 16$ and $M = 32$; for specific applications it may be beneficial to vary M along this tradeoff between quality and model size. This experiment uses the static *Lego* scene with 4 scales: 64, 128, 256, and 512.

Time smoothness regularizer. Sec. 3.2 describes our temporal smoothness regularizer based on penalizing the norm of the second derivative over the time dimension, to encourage smooth motion and discourage acceleration. Tab. 6 illustrates an ablation study of this regularizer on the *Jumping Jacks* scene from D-NeRF [30].

Time Smoothness Weight (λ)	Explicit PSNR \uparrow	Hybrid PSNR \uparrow
0.0	29.69	29.86
0.001	30.78	31.24
0.01	31.17	31.56
0.1	31.10	31.50
1.0	30.67	31.22
10.0	29.72	30.57

Table 6. **Ablation study over temporal smoothness regularization.** For both models, a temporal smoothness weight of 0.01 to 0.1 is best, with PSNR degrading with over- or under-regularization. This experiment uses the *Jumping Jacks* scene with 4 scales: 64, 128, 256, and 512, and 32 features per scale.

7. Model hyperparameters

Our full hyperparameter settings are available in the config files in our released code, at github.com/sarafridov/K-Planes.

PSNR \uparrow									
	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship	Mean
Ours-explicit	34.15	27.99	31.71	37.81	36.40	32.11	34.14	30.71	33.13
Ours-hybrid	34.74	27.84	32.40	38.17	37.21	32.66	34.76	31.19	33.62
INGP [25]	35.00	26.02	33.51	37.40	36.39	29.78	36.22	31.10	33.18
TensorRF [6]	35.76	26.01	33.99	37.41	36.46	30.12	34.61	30.77	33.14
Plenoxels [10]	33.98	25.35	31.83	36.43	34.10	29.14	33.26	29.62	31.71
JAXNeRF [7, 24]	34.20	25.27	31.15	36.81	34.02	30.30	33.72	29.33	31.85

SSIM \uparrow									
	Chair	Drums	Ficus	Hotdog	Lego	Materials	Mic	Ship	Mean
Ours-explicit	0.976	0.953	0.977	0.981	0.980	0.968	0.985	0.892	0.964
Ours-hybrid	0.980	0.953	0.980	0.982	0.983	0.971	0.987	0.899	0.967
INGP	-	-	-	-	-	-	-	-	-
TensorRF	0.985	0.937	0.982	0.982	0.983	0.952	0.988	0.895	0.963
Plenoxels	0.977	0.933	0.976	0.980	0.975	0.949	0.985	0.890	0.958
JAXNeRF	0.975	0.929	0.970	0.978	0.970	0.955	0.983	0.868	0.954

Table 7. **Full results on static synthetic scenes [24]**. Dashes denote values that were not reported in prior work.

PSNR \uparrow									
	Room	Fern	Leaves	Fortress	Orchids	Flower	T-Rex	Horns	Mean
Ours-explicit	32.72	24.87	21.07	31.34	19.89	28.37	27.54	28.40	26.78
Ours-hybrid	32.64	25.38	21.30	30.44	20.26	28.67	28.01	28.64	26.92
NeRF [24]	32.70	25.17	20.92	31.16	20.36	27.40	26.80	27.45	26.50
Plenoxels [10]	30.22	25.46	21.41	31.09	20.24	27.83	26.48	27.58	26.29
TensorRF (L) [6]	32.35	25.27	21.30	31.36	19.87	28.60	26.97	28.14	26.73
DVGOv2 [33]	-	-	-	-	-	-	-	-	26.34

SSIM \uparrow									
	Room	Fern	Leaves	Fortress	Orchids	Flower	T-Rex	Horns	Mean
Ours-explicit	0.955	0.809	0.738	0.898	0.665	0.867	0.909	0.884	0.841
Ours-hybrid	0.957	0.828	0.746	0.890	0.676	0.872	0.915	0.892	0.847
NeRF [24]	0.948	0.792	0.690	0.881	0.641	0.827	0.880	0.828	0.811
Plenoxels [10]	0.937	0.832	0.760	0.885	0.687	0.862	0.890	0.857	0.839
TensorRF (L) [6]	0.952	0.814	0.752	0.897	0.649	0.871	0.900	0.877	0.839
DVGOv2 [33]	-	-	-	-	-	-	-	-	0.838

Table 8. **Full results on static forward-facing scenes [23]**. Dashes denote values that were not reported in prior work.

PSNR \uparrow									
	Hell Warrior	Mutant	Hook	Balls	Lego	T-Rex	Stand Up	Jumping Jacks	Mean
Ours-explicit	24.77	32.43	27.84	38.59	25.25	30.45	32.50	31.26	30.39
Ours-hybrid	24.81	32.59	28.13	40.33	25.27	30.75	33.17	31.64	30.84
D-NeRF [30]	25.02	31.29	29.25	32.80	21.64	31.75	32.79	32.80	29.67
T-NeRF [30]	23.19	30.56	27.21	32.01	23.82	30.19	31.24	32.01	28.78
Tensor4D [†] [31]	-	-	-	-	26.71	-	36.32	34.43	-
TiNeuVox [†] [9]	28.17	33.61	31.45	40.73	25.02	32.70	35.43	34.23	32.67
V4D [†] [11]	27.03	36.27	31.04	42.67	25.62	34.53	37.20	35.36	33.72

SSIM \uparrow									
	Hell Warrior	Mutant	Hook	Balls	Lego	T-Rex	Stand Up	Jumping Jacks	Mean
Ours-explicit	0.948	0.968	0.941	0.986	0.934	0.972	0.973	0.966	0.961
Ours-hybrid	0.950	0.968	0.945	0.989	0.936	0.973	0.977	0.969	0.963
D-NeRF [30]	0.95	0.97	0.96	0.98	0.83	0.97	0.98	0.98	0.95
T-NeRF [30]	0.93	0.96	0.94	0.97	0.90	0.96	0.97	0.97	0.95
Tensor4D [31]	-	-	-	-	0.953	-	0.983	0.982	-
TiNeuVox [9]	0.97	0.98	0.97	0.99	0.92	0.98	0.99	0.98	0.97
V4D [11]	0.96	0.99	0.97	0.99	0.95	0.99	0.99	0.99	0.98

[†] Very recent/concurrent work. Tensor4D was released in November 2022, TiNeuVox in September 2022, and V4D in October 2022. Note that TiNeuVox trains and evaluates using half-resolution, and Tensor4D and V4D have not yet released code, so their train/test resolution is unknown.

Table 9. **Full results on monocular “teleporting-camera” dynamic scenes.** We use the synthetic scenes from D-NeRF [30], which we refer to as monocular “teleporting-camera” because although there is a single training view per timestep, the camera can move arbitrarily between adjacent timesteps. Dashes denote unreported values. TiNeuVox trains in 30 minutes, V4D in 4.9 hours, D-NeRF in 2 days, and Tensor4D for an unspecified duration (Tensor4D reports iterations rather than time). Our reported results were obtained after roughly 1 hour of optimization on a single GPU.

PSNR \uparrow							
	Coffee Martini	Spinach	Cut Beef	Flame Salmon ¹	Flame Steak	Sear Steak	Mean
Ours-explicit	28.74	32.19	31.93	28.71	31.80	31.89	30.88
Ours-hybrid	29.99	32.60	31.82	30.44	32.38	32.52	31.63
LLFF [23]	-	-	-	23.24	-	-	-
DyNeRF [16]	-	-	-	29.58	-	-	-
MixVoxels-L [†] [37]	29.36	31.61	31.30	29.92	31.21	31.43	30.80

SSIM \uparrow							
	Coffee Martini	Cook Spinach	Cut Beef	Flame Salmon ¹	Flame Steak	Sear Steak	Mean
Ours-explicit	0.943	0.968	0.965	0.942	0.970	0.971	0.960
Ours-hybrid	0.953	0.966	0.966	0.953	0.970	0.974	0.964
LLFF	-	-	-	0.848	-	-	-
DyNeRF	-	-	-	0.961	-	-	-
MixVoxels-L	0.946	0.965	0.965	0.945	0.970	0.971	0.960

[†] Very recent/concurrent work. MixVoxels was released in December 2022. ¹Using the first 10 seconds of the 30 second long video.

Table 10. **Full results on multiview dynamic scenes [16].** Dashes denote unreported values. Note that our method optimizes in less than 4 GPU hours, whereas DyNeRF trains on 8 GPUs for a week, approximately 1344 GPU hours.

PSNR \uparrow				
	Brandenburg Gate	Sacre Coeur	Trevi Fountain	Mean
Ours-explicit	24.85	19.90	22.00	22.25
Ours-hybrid	25.49	20.61	22.67	22.92
NeRF-W [20]	29.08	25.34	26.58	27.00
NeRF-W (public) [†]	21.32	19.17	18.61	19.70
LearnIt [35]	19.11	19.33	19.35	19.26
MS-SSIM \uparrow				
	Brandenburg Gate	Sacre Coeur	Trevi Fountain	Mean
Ours-explicit	0.912	0.821	0.845	0.859
Ours-hybrid	0.924	0.852	0.856	0.877
NeRF-W	0.962	0.939	0.934	0.945
NeRF-W (public) [†]	0.845	0.752	0.694	0.764
LearnIt	-	-	-	-

[†] Open-source version https://github.com/kweal23/nerf_pl/tree/nerfw where we implement the test-time optimization ourselves exactly as for k -planes. NeRF-W code is not public.

Table 11. **Full results on phototourism scenes.** Note that our results were obtained after about 35 GPU *minutes*, whereas NeRF-W trains with 8 GPUs for two days, approximately 384 GPU hours.