

TensorRF: Tensorial Radiance Fields

Anpei Chen^{1*} Zexiang Xu^{2*} Andreas Geiger³ Jingyi Yu¹ Hao Su⁴

¹ShanghaiTech University ²Adobe Research

³University of Tübingen and MPI-IS, Tübingen ⁴UC San Diego

<https://apchenstu.github.io/TensorRF/>

Abstract. We present TensorRF, a novel approach to model and reconstruct radiance fields. Unlike NeRF that purely uses MLPs, we model the radiance field of a scene as a 4D tensor, which represents a 3D voxel grid with per-voxel multi-channel features. Our central idea is to factorize the 4D scene tensor into multiple compact low-rank tensor components. We demonstrate that applying traditional CP decomposition – that factorizes tensors into rank-one components with compact vectors – in our framework leads to improvements over vanilla NeRF. To further boost performance, we introduce a novel vector-matrix (VM) decomposition that relaxes the low-rank constraints for two modes of a tensor and factorizes tensors into compact vector and matrix factors. Beyond superior rendering quality, our models with CP and VM decompositions lead to a significantly lower memory footprint in comparison to previous and concurrent works that directly optimize per-voxel features. Experimentally, we demonstrate that TensorRF with CP decomposition achieves fast reconstruction (< 30 min) with better rendering quality and even a smaller model size (< 4 MB) compared to NeRF. Moreover, TensorRF with VM decomposition further boosts rendering quality and outperforms previous state-of-the-art methods, while reducing the reconstruction time (< 10 min) and retaining a compact model size (< 75 MB).

1 Introduction

Modeling and reconstructing 3D scenes as representations that support high-quality image synthesis is crucial for computer vision and graphics with various applications in visual effects, e-commerce, virtual and augmented reality, and robotics. Recently, NeRF [24] and its many follow-up works [49,18] have shown success on modeling a scene as a radiance field and enabled photo-realistic rendering of scenes with highly complex geometry and view-dependent appearance effects. Despite the fact that (purely MLP-based) NeRF models require small memory, they take a long time (hours or days) to train. In this work, we pursue a novel approach that is both *efficient in training time* and *compact in memory footprint*, and at the same time achieves *state-of-the-art* rendering quality.

* Equal contribution.

Research done when Anpei Chen was in a remote internship with UCSD.

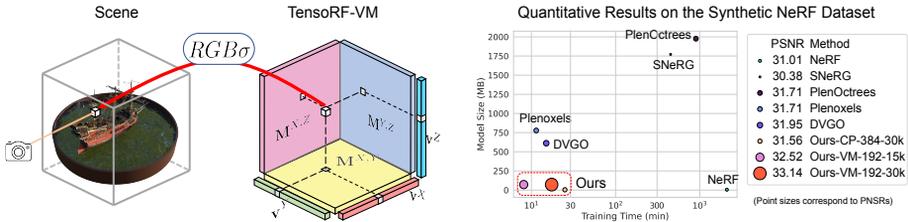


Fig. 1: Left: We model a scene as a tensorial radiance field using a set of vectors (\mathbf{v}) and matrices (\mathbf{M}) that describe scene appearance and geometry along their corresponding axes. These vector/matrix factors are used to compute volume density σ and view-dependent RGB color via vector-matrix outer products for realistic volume rendering. Right: In comparison with previous and concurrent methods, our TensorRF models can achieve the best rendering quality and are the only methods that can simultaneously achieve fast reconstruction and high compactness. (Our models are denoted with their decomposition techniques, number of components, and training steps.)

To do so, we propose TensorRF, a novel radiance field representation that is highly compact and also fast to reconstruct, enabling efficient scene reconstruction and modeling. Unlike coordinate-based MLPs used in NeRF, we represent radiance fields as an explicit voxel grid of features. Note that it is *unclear* whether voxel grid representation can benefit the efficiency of reconstruction: While previous work has used feature grids [18,47,12], they require large GPU memory to store the voxels whose size grows cubically with resolution, and some even require pre-computing a NeRF for distillation, leading to very long reconstruction time.

Our work addresses the inefficiency of voxel grid representations in a principled framework, leading to a family of simple yet effective methods. We leverage the fact that a feature grid can naturally be seen as a 4D tensor, where three of its modes correspond to the XYZ axes of the grid and the fourth mode represents the feature channel dimension. This opens the possibility of exploiting classical tensor decomposition techniques – which have been widely applied to high-dimensional data analysis and compression in various fields [16] – for radiance field modeling. We, therefore, propose to factorize the tensor of radiance fields into multiple *low-rank* tensor components, leading to an accurate and compact scene representation. Note that our central idea of tensorizing radiance fields is general and can be potentially adopted to any tensor decomposition technique.

In this work, we first attempt the classic CANDECOMP/PARAFAC (CP) decomposition [5]. We show that TensorRF with CP decomposition can already achieve photo-realistic rendering and lead to a more compact model than NeRF that is purely MLP based (see Fig. 1 and Tab. 1). However, experimentally, to further push reconstruction quality for complex scenes, we have to use more component factors, which undesirably increases training time.

Therefore, we present a novel vector-matrix (VM) decomposition technique that effectively reduces the number of components required for the same expression capacity, leading to faster reconstruction and better rendering. In particular,

inspired by the CP and block term decomposition [9], we propose to factorize the full tensor of a radiance field into multiple vector and matrix factors per tensor component. Unlike the sum of outer products of pure vectors in CP decomposition, we consider the sum of vector-matrix outer products (see Fig. 2). In essence, we relax the ranks of two modes of each component by jointly modeling two modes in a matrix factor. While this increases the model size compared to pure vector-based factorization in CP, we enable each component to express more complex tensor data of higher ranks, thus significantly reducing the required number of components in radiance field modeling and leading to better rendering quality (see Fig. 1 and Tab. 2).

With CP/VM decomposition, our approach compactly encodes spatially varying features in the voxel grid. Volume density and view-dependent color can be decoded from the features, supporting volumetric radiance field rendering. Because a tensor expresses discrete data, we also enable efficient trilinear interpolation for our representation to model a continuous field. Our representation supports various types of per-voxel features with different decoding functions, including neural features – depending on an MLP to regress view-dependent colors from the features – and spherical harmonics (SH) features (coefficients) – allowing for simple color computation from the fixed SH functions and leading to a representation without neural networks.

Our tensorial radiance fields can be effectively reconstructed from multi-view images and enable realistic novel view synthesis. In contrast to previous works that directly reconstruct voxels, our tensor factorization reduces space complexity from $\mathcal{O}(n^3)$ to $\mathcal{O}(n)$ (with CP) or $\mathcal{O}(n^2)$ (with VM), significantly lowering memory footprint. We present extensive evaluation of our approach with various settings, covering both CP and VM models, different numbers of components and grid resolutions. We demonstrate that all models are able to achieve realistic novel view synthesis results that are on par or better than previous state-of-the-art methods (see Fig. 1 and Tab. 1). More importantly, our approach is of high computation and memory efficiency. All TensorRF models can reconstruct high-quality radiance fields in 30 min; our fastest model with VM decomposition takes less than 10 min, which is significantly faster (about 100x) than NeRF and many other methods, while requiring substantially less memory than previous and concurrent voxel-based methods. Note that, unlike concurrent works [46,25] that require unique data structures and customized CUDA kernels, our model’s efficiency gains are obtained using a standard PyTorch implementation. As far as we know, our work is the first that views radiance field modeling from a tensorial perspective and pose the problem of radiance field reconstruction as one of low-rank tensor reconstructions.

2 Related Work

Tensor decomposition. Tensor decomposition [16] has been studied for decades in various fields. In general, the most widely used decompositions are Tucker decomposition [38] and CP decomposition [5,11], both of which can be seen as

generalizations of the matrix singular value decomposition (SVD). CP decomposition can also be seen as a special Tucker decomposition whose core tensor is diagonal. By combining CP and Tucker decomposition, block term decomposition (BTD) has been proposed with its many variants [9] and used in many vision and learning tasks [45,1,44]. In this work, we leverage tensor decomposition for radiance field modeling. We directly apply CP decomposition and also introduce a new vector-matrix decomposition, which can be seen as a special BTD.

Scene representations and radiance fields. Various scene representations, including meshes [10,39], point clouds [31], volumes [13,32], implicit functions [22,30], have been extensively studied in recent years. Many neural representations [50,35,20,3] are proposed for high-quality rendering. NeRF [24] introduces radiance fields to address novel view synthesis and achieves photo-realistic quality. This representation has been quickly extended and applied in diverse graphics and vision applications, including generative models [6,26], appearance acquisition [2,4], surface reconstruction [40,27], fast rendering [33,47,12], appearance editing [42,19], dynamic capture [17,28]. While leading to realistic rendering and a compact model, NeRF with its pure MLP-based representation has known limitations in slow reconstruction and rendering. Recent methods [47,18,12] have leveraged a voxel grid of features in radiance field modeling, achieving fast rendering. However, these grid-based methods still require long reconstruction time and even lead to high memory costs, sacrificing the compactness of NeRF. Based on feature grids, we present a novel tensorial scene representation, leveraging tensor factorization techniques, leading to fast reconstruction and compact modeling.

Other methods design generalizable network modules trained across scenes to achieve image-dependent radiance field rendering [48,41,8] and fast reconstruction [7,43]. Our approach focuses on radiance field representation and only considers per-scene optimization (like NeRF). We show that our representation can already lead to highly efficient radiance field reconstruction without any across-scene generalization. We leave the extensions to generalizable settings as future work.

Concurrent work. The field of radiance field modeling is moving very fast and many concurrent works have appeared on arXiv as preprints over the last three months. DVGO [37] and Plenoxels [46] also optimize voxel grids of (neural or SH) features for fast radiance field reconstruction. However, they still optimize per-voxel features directly like previous voxel-based methods, thus requiring large memory. Our approach instead factorizes the feature grid into compact components and leads to significantly higher memory efficiency. Instant-NGP [25] uses multi-resolution hashing for efficient encoding and also leads to high compactness. This technique is orthogonal to our factorization-based technique; potentially, each of our vector/matrix factor can be encoded with this hashing technique and we leave such combination as future work.

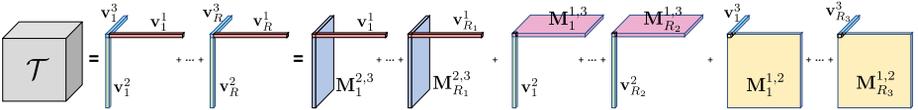


Fig. 2: Tensor factorization. Left: CP decomposition (Eqn. 1), which factorizes a tensor as a sum of vector outer products. Right: our vector-matrix decomposition (Eqn. 3), which factorizes a tensor as a sum of vector-matrix outer products.

3 CP and VM Decomposition

We factorize radiance fields into compact components for scene modeling. To do so, we apply both the classic CP decomposition and a new vector-matrix (VM) decomposition; both are illustrated in Fig. 2. We now discuss both decompositions with an example of a 3D (3rd-order) tensor. We will introduce how to apply tensor factorization in radiance field modeling (with a 4D tensor) in Sec. 4.

CP decomposition. Given a 3D tensor $\mathcal{T} \in \mathbb{R}^{I \times J \times K}$, CP decomposition factorizes it into a sum of outer products of vectors (shown in Fig. 2):

$$\mathcal{T} = \sum_{r=1}^R \mathbf{v}_r^1 \circ \mathbf{v}_r^2 \circ \mathbf{v}_r^3 \quad (1)$$

where $\mathbf{v}_r^1 \circ \mathbf{v}_r^2 \circ \mathbf{v}_r^3$ corresponds to a rank-one tensor component, and $\mathbf{v}_r^1 \in \mathbb{R}^I$, $\mathbf{v}_r^2 \in \mathbb{R}^J$, and $\mathbf{v}_r^3 \in \mathbb{R}^K$ are factorized vectors of the three modes for the r th component. Superscripts denote the modes of each factor; \circ represents the outer product. Hence, each tensor element \mathcal{T}_{ijk} is a sum of scalar products:

$$\mathcal{T}_{ijk} = \sum_{r=1}^R \mathbf{v}_{r,i}^1 \mathbf{v}_{r,j}^2 \mathbf{v}_{r,k}^3 \quad (2)$$

where i, j, k denote the indices of the three modes.

CP decomposition factorizes a tensor into multiple vectors, expressing multiple compact rank-one components. CP can be directly applied in our tensorial radiance field modeling and generate high-quality results (see Tab. 1). However, because of too high compactness, CP decomposition can require many components to model complex scenes, leading to high computational costs in radiance field reconstruction. Inspired by block term decomposition (BTD), we present a new VM decomposition, leading to more efficient radiance field reconstruction.

Vector-Matrix (VM) decomposition. Unlike CP decomposition that utilizes pure vector factors, VM decomposition factorizes a tensor into multiple vectors and matrices as shown in Fig. 2 right. This is expressed by

$$\mathcal{T} = \sum_{r=1}^{R_1} \mathbf{v}_r^1 \circ \mathbf{M}_r^{2,3} + \sum_{r=1}^{R_2} \mathbf{v}_r^2 \circ \mathbf{M}_r^{1,3} + \sum_{r=1}^{R_3} \mathbf{v}_r^3 \circ \mathbf{M}_r^{1,2} \quad (3)$$

where $\mathbf{M}_r^{2,3} \in \mathbb{R}^{J \times K}$, $\mathbf{M}_r^{1,3} \in \mathbb{R}^{I \times K}$, $\mathbf{M}_r^{1,2} \in \mathbb{R}^{I \times J}$ are matrix factors for two (denoted by superscripts) of the three modes. For each component, we relax

its two mode ranks to be arbitrarily large, while restricting the third mode to be rank-one; e.g., for component tensor $\mathbf{v}_r^1 \circ \mathbf{M}_r^{2,3}$, its mode-1 rank is 1, and its mode-2 and mode-3 ranks can be arbitrary, depending on the rank of the matrix $\mathbf{M}_r^{2,3}$. In general, instead of using separate vectors in CP, we combine every two modes and represent them by matrices, allowing each mode to be adequately parametrized with a smaller number of components. R_1, R_2, R_3 can be set differently and should be chosen depending on the complexity of each mode. Our VM decomposition can be seen as a special case of general BTD.

Note that, each of our component tensors has more parameters than a component in CP decomposition. While this leads to lower compactness, a VM component tensor can express more complex high-dimensional data than a CP component, thus reducing the required number of components when modeling the same complex function. On the other hand, VM decomposition is still of very high compactness, reducing memory complexity from $\mathcal{O}(N^3)$ to $\mathcal{O}(N^2)$, compared to dense grid representations.

Tensor for scene modeling. In this work, we focus on the task of modeling and reconstructing radiance fields. In this case, the three tensor modes correspond to XYZ axes, and we thus directly denote the modes with XYZ to make it intuitive. Meanwhile, in the context of 3D scene representation, we consider $R_1 = R_2 = R_3 = R$ for most of scenes, reflecting the fact that a scene can distribute and appear equally complex along its three axes. Therefore, Eqn. 3 can be re-written as

$$\mathcal{T} = \sum_{r=1}^R \mathbf{v}_r^X \circ \mathbf{M}_r^{Y,Z} + \mathbf{v}_r^Y \circ \mathbf{M}_r^{X,Z} + \mathbf{v}_r^Z \circ \mathbf{M}_r^{X,Y} \quad (4)$$

In addition, to simplify notation and the following discussion in later sections, we also denote the three types of component tensors as $\mathcal{A}_r^X = \mathbf{v}_r^X \circ \mathbf{M}_r^{Y,Z}$, $\mathcal{A}_r^Y = \mathbf{v}_r^Y \circ \mathbf{M}_r^{X,Z}$, and $\mathcal{A}_r^Z = \mathbf{v}_r^Z \circ \mathbf{M}_r^{X,Y}$; here the superscripts XYZ of \mathcal{A} indicate different types of components. With this, a tensor element \mathcal{T}_{ijk} is expressed as

$$\mathcal{T}_{ijk} = \sum_{r=1}^R \sum_m \mathcal{A}_{r,ijk}^m \quad (5)$$

where $m \in XYZ$, $\mathcal{A}_{r,ijk}^X = \mathbf{v}_{r,i}^X \mathbf{M}_{r,jk}^{YZ}$, $\mathcal{A}_{r,ijk}^Y = \mathbf{v}_{r,j}^Y \mathbf{M}_{r,ik}^{XZ}$, and $\mathcal{A}_{r,ijk}^Z = \mathbf{v}_{r,k}^Z \mathbf{M}_{r,ij}^{XY}$. Similarly, we can also denote a CP component as $\mathcal{A}^\gamma = \mathbf{v}_r^X \circ \mathbf{v}_r^Y \circ \mathbf{v}_r^Z$, and Eqn. 5 can also express CP decomposition by considering $m = \gamma$, where the summation over m can be removed.

4 Tensorial Radiance Field Representation

We now present our Tensorial Radiance Field Representation (TensorRF). For simplicity, we focus on presenting TensorRF with our VM decomposition. CP decomposition is simpler and its decomposition equations can be directly applied with minimal modification (like Eqn. 5).

4.1 Feature grids and radiance field.

Our goal is to model a radiance field, which is essentially a function that maps any 3D location \mathbf{x} and viewing direction d to its volume density σ and view-dependent color c , supporting differentiable ray marching for volume rendering. We leverage a regular 3D grid \mathcal{G} with per-voxel multi-channel features to model such a function. We split it (by feature channels) into a geometry grid \mathcal{G}_σ and an appearance grid \mathcal{G}_c , separately modelling the volume density σ and view-dependent color c .

Our approach supports various types of appearance features in \mathcal{G}_c , depending on a pre-selected function S that converts an appearance feature vector and a viewing direction d to color c . For example, S can be a small MLP or spherical harmonics (SH) functions, where \mathcal{G}_c contains neural features and SH coefficients respectively. We show that both MLP and SH functions work well with our model (see Tab.1). On the other hand, we consider a single-channel grid \mathcal{G}_σ , whose values represent volume density directly, without requiring an extra converting function. The continuous grid-based radiance field can be written as

$$\sigma, c = \mathcal{G}_\sigma(\mathbf{x}), S(\mathcal{G}_c(\mathbf{x}), d) \quad (6)$$

where $\mathcal{G}_\sigma(\mathbf{x})$, $\mathcal{G}_c(\mathbf{x})$ represent the trilinearly interpolated features from the two grids at location \mathbf{x} . We model \mathcal{G}_σ and \mathcal{G}_c as factorized tensors.

4.2 Factorizing radiance fields

While $\mathcal{G}_\sigma \in \mathbb{R}^{I \times J \times K}$ is a 3D tensor, $\mathcal{G}_c \in \mathbb{R}^{I \times J \times K \times P}$ is a 4D tensor. Here I, J, K correspond to the resolutions of the feature grid along the X, Y, Z axes, and P is the number of appearance feature channels.

We factorize these radiance field tensors to compact components. In particular, with the VM decomposition, the 3D geometry tensor \mathcal{G}_σ is factorized as

$$\mathcal{G}_\sigma = \sum_{r=1}^{R_\sigma} \mathbf{v}_{\sigma,r}^X \circ \mathbf{M}_{\sigma,r}^{YZ} + \mathbf{v}_{\sigma,r}^Y \circ \mathbf{M}_{\sigma,r}^{XZ} + \mathbf{v}_{\sigma,r}^Z \circ \mathbf{M}_{\sigma,r}^{XY} = \sum_{r=1}^{R_\sigma} \sum_{m \in XYZ} \mathcal{A}_{\sigma,r}^m \quad (7)$$

The appearance tensor \mathcal{G}_c has an additional mode corresponding to the feature channel dimension. Note that, compared to the XYZ modes, this mode is often of lower dimension, leading to a lower rank. Therefore, we do not combine this mode with other modes in matrix factors and instead only use vectors, denoted by \mathbf{b}_r , for this mode in the factorization. Specifically, \mathcal{G}_c is factorized as

$$\begin{aligned} \mathcal{G}_c &= \sum_{r=1}^{R_c} \mathbf{v}_{c,r}^X \circ \mathbf{M}_{c,r}^{YZ} \circ \mathbf{b}_{3r-2} + \mathbf{v}_{c,r}^Y \circ \mathbf{M}_{c,r}^{XZ} \circ \mathbf{b}_{3r-1} + \mathbf{v}_{c,r}^Z \circ \mathbf{M}_{c,r}^{XY} \circ \mathbf{b}_{3r} \\ &= \sum_{r=1}^{R_c} \mathcal{A}_{c,r}^X \circ \mathbf{b}_{3r-2} + \mathcal{A}_{c,r}^Y \circ \mathbf{b}_{3r-1} + \mathcal{A}_{c,r}^Z \circ \mathbf{b}_{3r} \end{aligned} \quad (8)$$

Note that, we have $3R_c$ vectors \mathbf{b}_r to match the total number of components.

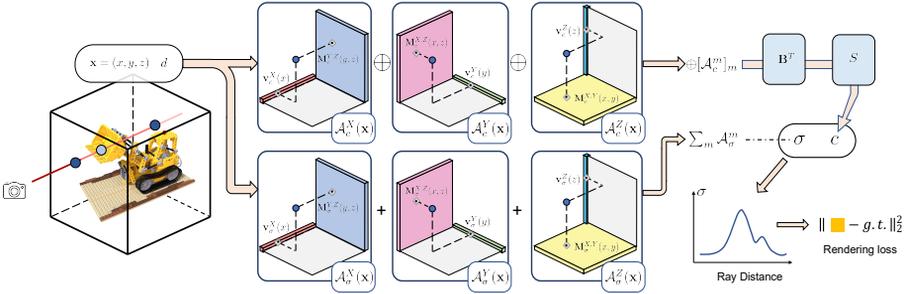


Fig. 3: TensorRF (VM) reconstruction and rendering. We model radiance fields as tensors using a set of vectors (\mathbf{v}) and matrices (\mathbf{M}), which describe the scene along their corresponding (XYZ) axes and are used for computing volume density σ and view-dependent color c in differentiable ray marching. For each shading location $\mathbf{x} = (x, y, z)$, we use linearly/bilinearly sampled values from the vector/matrix factors to efficiently compute the corresponding trilinearly interpolated values ($\mathcal{A}(\mathbf{x})$) of the tensor components. The density component values ($\mathcal{A}_\sigma(\mathbf{x})$) are summed to get the volume density (σ) directly. The appearance values ($\mathcal{A}_c(\mathbf{x})$) are concatenated into a vector ($\oplus[\mathcal{A}_c^m(\mathbf{x})]_m$) that is then multiplied by an appearance matrix \mathbf{B} and sent to the decoding function S for RGB color (c) regression.

Overall, we factorize the entire tensorial radiance field into $3R_\sigma + 3R_c$ matrices ($\mathbf{M}_{\sigma,r}^{YZ}, \dots, \mathbf{M}_{c,r}^{YZ}, \dots$) and $3R_\sigma + 6R_c$ vectors ($\mathbf{v}_{\sigma,r}^X, \dots, \mathbf{v}_{c,r}^X, \dots, \mathbf{b}_r$). In general, we adopt $R_\sigma \ll I, J, K$ and $R_c \ll I, J, K$, leading to a highly compact representation that can encode a high-resolution dense grid. In essence, the XYZ-mode vector and matrix factors, $\mathbf{v}_{\sigma,r}^X$, $\mathbf{M}_{\sigma,r}^{YZ}$, $\mathbf{v}_{c,r}^X$, $\mathbf{M}_{c,r}^{YZ}$, ..., describe the spatial distributions of the scene geometry and appearance along their corresponding axes. On the other hand, the appearance feature-mode vectors \mathbf{b}_r express the global appearance correlations. By stacking all \mathbf{b}_r as columns together, we have a $P \times 3R_c$ matrix \mathbf{B} ; this matrix \mathbf{B} can also be seen as a global appearance dictionary that abstracts the appearance commonalities across the entire scene.

4.3 Efficient feature evaluation.

Our factorization-based model can compute each voxel's feature vector at low costs, only requiring one value per XYZ-mode vector/matrix factor. We also enable efficient trilinear interpolation for our model, leading to a continuous field.

Direct evaluation. With VM factorization, a density value $\mathcal{G}_{\sigma,ijk}$ of a single voxel at indices ijk can be directly and efficiently evaluated by following Eqn. 5:

$$\mathcal{G}_{\sigma,ijk} = \sum_{r=1}^{R_\sigma} \sum_{m \in XYZ} \mathcal{A}_{\sigma,r,ijk}^m \quad (9)$$

Here, computing each $\mathcal{A}_{\sigma,r,ijk}^m$ only requires indexing and multiplying two values from its corresponding vector and matrix factors.

As for the appearance grid \mathcal{G}_c , we always need to compute a full P -channel feature vector, which the shading function S requires as input, corresponding to a 1D slice of \mathcal{G}_c at fixed XYZ indices ijk :

$$\mathcal{G}_{c,ijk} = \sum_{r=1}^{R_c} \mathcal{A}_{c,r,ijk}^X \mathbf{b}_{3r-2} + \mathcal{A}_{c,r,ijk}^Y \mathbf{b}_{3r-1} + \mathcal{A}_{c,r,ijk}^Z \mathbf{b}_{3r} \quad (10)$$

Here, there's no additional indexing for the feature mode, since we compute a full vector. We further simplify Eqn. 10 by re-ordering the computation. For this, we denote $\oplus[\mathcal{A}_{c,ijk}^m]_{m,r}$ as the vector that stacks all $\mathcal{A}_{c,r,ijk}^m$ values for $m = X, Y, Z$ and $r = 1, \dots, R_c$, which is a vector of $3R_c$ dimensions; \oplus can also be considered as the concatenation operator that concatenates all scalar values (1-channel vectors) into a $3R_c$ -channel vector in practice. Using matrix \mathbf{B} (introduced in Sec. 4.1) that stacks all \mathbf{b}_r , Eqn. 10 is equivalent to a matrix vector product:

$$\mathcal{G}_{c,ijk} = \mathbf{B}(\oplus[\mathcal{A}_{c,ijk}^m]_{m,r}) \quad (11)$$

Note that, Eqn. 11 is not only formally simpler but also leads to a simpler implementation in practice. Specifically, when computing a large number of voxels in parallel, we first compute and concatenate $\mathcal{A}_{c,r,ijk}^m$ for all voxels as column vectors in a matrix and then multiply the shared matrix \mathbf{B} once.

Trilinear interpolation. We apply trilinear interpolation to model a continuous field. Naïvely achieving trilinear interpolation is costly, as it requires evaluation of 8 tensor values and interpolating them, increasing computation by a factor of 8 compared to computing a single tensor element. However, we find that trilinearly interpolating a component tensor is naturally equivalent to interpolating its vector/matrix factors linearly/bilinearly for the corresponding modes, thanks to the beauty of linearity of the trilinear interpolation and the outer product.

For example, given a component tensor $\mathcal{A}_r^X = \mathbf{v}_r^X \circ \mathbf{M}_r^{YZ}$ with its each tensor element $\mathcal{A}_{r,ijk}^X = \mathbf{v}_{r,i}^X \mathbf{M}_{r,jk}^{YZ}$, we can compute its interpolated values as:

$$\mathcal{A}_r^X(\mathbf{x}) = \mathbf{v}_r^X(x) \mathbf{M}_r^{YZ}(y, z) \quad (12)$$

where $\mathcal{A}_r^X(\mathbf{x})$ is \mathcal{A}_r 's trilinearly interpolated value at location $\mathbf{x} = (x, y, z)$ in the 3D space, $\mathbf{v}_r^X(x)$ is \mathbf{v}_r^X 's linearly interpolated value at x along X axis, and $\mathbf{M}_r^{YZ}(y, z)$ is \mathbf{M}_r^{YZ} 's bilinearly interpolated value at (y, z) in the YZ plane. Similarly, we have $\mathcal{A}_r^Y(\mathbf{x}) = \mathbf{v}_r^Y(y) \mathbf{M}_r^{XZ}(x, z)$ and $\mathcal{A}_r^Z(\mathbf{x}) = \mathbf{v}_r^Z(z) \mathbf{M}_r^{XY}(x, y)$ (for CP decomposition $\mathcal{A}_r^Y(\mathbf{x}) = \mathbf{v}_r^X(x) \mathbf{v}_r^Y(y) \mathbf{v}_r^Z(z)$ is also valid). Thus, trilinearly interpolating the two grids is expressed as

$$\mathcal{G}_\sigma(\mathbf{x}) = \sum_r \sum_m \mathcal{A}_{\sigma,r}^m(\mathbf{x}) \quad (13)$$

$$\mathcal{G}_c(\mathbf{x}) = \mathbf{B}(\oplus[\mathcal{A}_{c,r}^m(\mathbf{x})]_{m,r}) \quad (14)$$

These equations are very similar to Eqn. 9 and 11, simply replacing the tensor elements with interpolated values. We avoid recovering 8 individual tensor elements for trilinear interpolation and instead directly recover the interpolated value, leading to low computation and memory costs at run time.

4.4 Rendering and reconstruction.

Equations 6, 12–14 describe the core components of our model. By combining Eqn. 6,13,14, our factorized tensorial radiance field can be expressed as

$$\sigma, c = \sum_r \sum_m \mathcal{A}_{\sigma,r}^m(\mathbf{x}), S(\mathbf{B}(\oplus[\mathcal{A}_{c,r}^m(\mathbf{x})]_{m,r}), d) \quad (15)$$

i.e., we obtain continuous volume density and view-dependent color given any 3D location and viewing direction. This allows for high-quality radiance field reconstruction and rendering. Note that, this equation is general and describes TensorRF with both CP and VM decomposition. Our full pipeline of radiance field reconstruction and rendering with VM decomposition is illustrated in Fig. 3.

Volume rendering. To render images, we use differentiable volume rendering, following NeRF [24]. Specifically, for each pixel, we march along a ray, sampling Q shading points along the ray and computing the pixel color by

$$C = \sum_{q=1}^Q \tau_q (1 - \exp(-\sigma_q \Delta_q)) c_q, \quad \tau_q = \exp(-\sum_{p=1}^{q-1} \sigma_p \Delta_p) \quad (16)$$

Here, σ_q, c_q are the corresponding density and color computed by our model at their sampled locations \mathbf{x}_q ; Δ_q is the ray step size and τ_q represents transmittance.

Reconstruction. Given a set of multi-view input images with known camera poses, our tensorial radiance field is optimized per scene via gradient descent, minimizing an L2 rendering loss, using only the ground truth pixel colors as supervision. Our radiance field is explained by tensor factorization and modeled by a set of global vectors and matrices as basis factors that correlate and regularize the entire field in the optimization. However, this can sometimes lead to overfitting and local minima issues in gradient descent, leading to outliers or noises in regions with fewer observations. We utilize standard regularization terms that are commonly used in compressive sensing, including an L1 norm loss and a TV (total variation) loss on our vector and matrix factors, which effectively address these issues. We find that only applying the L1 sparsity loss is adequate for most datasets. However, for real datasets that have very few input images (like LLFF[23]) or imperfect capture conditions (like Tanks and Temples [15,18] that has varying exposure and inconsistent masks), a TV loss is more efficient than the L1 norm loss.

To further improve quality and avoid local minima, we apply coarse-to-fine reconstruction. Unlike previous coarse-to-fine techniques that require unique subdivisions on their sparse chosen sets of voxels, our coarse-to-fine reconstruction is simply achieved by linearly and bilinearly upsampling our XYZ-mode vector and matrix factors.

5 Implementation details

We briefly discuss our implementation; please refer to the appendix for more details. We implement our TensorRF using PyTorch [29], without customized

CUDA kernels. We implement the feature decoding function S as either an MLP or SH function and use $P = 27$ features for both. For SH, this corresponds to 3rd-order SH coefficients with RGB channels. For neural features, we use a small MLP with two FC layers (with 128-channel hidden layers) and ReLU activation.

We use the Adam optimizer [14] with initial learning rates of 0.02 for tensor factors and (when using neural features) 0.001 for the MLP decoder. We optimize our model for T steps with a batch size of 4096 pixel rays on a single Tesla V100 GPU (16GB). We apply a feature grid with a total number of N^3 voxels; the actual resolution of each dimension is computed based on the shape of the bounding box. To achieve coarse-to-fine reconstruction, we start from an initial low-resolution grid with N_0^3 voxels with $N_0 = 128$; we then upsample the vectors and matrices linearly and bilinearly at steps 2000, 3000, 4000, 5500, 7000 with the numbers of voxels interpolated between N_0^3 and N^3 linearly in logarithmic space. Please refer to Sec. 6 for the ablation study on different total steps (T), different resolutions (N), and different number of total components ($3R_\sigma + 3R_c$).

To facilitate reconstruction, we compute a binary occupancy mask grid at steps 2000 and 4000 using the volume density prediction from the intermediate TensorRF model to avoid computation in empty space. For datasets that do not provide bounding boxes, we start from a conservatively large box and leverage the occupancy mask computed at step 2000 to re-compute a more compact bounding box, with which we shrink and resample our tensor factors, leading to more precise modeling. For forward-facing scenes in the LLFF dataset [23], we apply NDC transformation that bounds the scene in a perspective frustum.

6 Experiments

We now present an extensive evaluation of our tensorial radiance fields. We first analyze our decomposition techniques, the number of components, grid resolutions, and optimization steps. We then compare our approach with previous and concurrent works on both 360° objects and forward-facing datasets.

Analysis of different TensorRF models. We evaluate our TensorRF on the Synthetic NeRF dataset [24] using both CP and VM decompositions with different numbers of components and different numbers of grid voxels. Table 2 shows the averaged rendering PSNRs, reconstruction time, and model size for each model. We use the same MLP decoding function (as described in Sec. 5) for all variants and optimize each model for 30k steps with a batch size of 4096.

Note that both TensorRF-CP and TensorRF-VM achieve consistently better rendering quality with more components or higher grid resolutions. TensorRF-CP achieves super compact modeling; even the largest model with 384 components and 500^3 voxels requires less than 4MB. This CP model also achieves the best rendering quality in all of our CP variants, leading to a high PSNR of 31.56, which even outperforms vanilla NeRF (see Tab. 1).

On the other hand, because it compresses more parameters in each component, TensorRF-VM achieves significantly better rendering quality than TensorRF-CP; even the smallest TensorRF-VM model with only 48 components and 200^3 voxels

Method	BatchSize	Steps	Synthetic-NeRF		NSVF		TanksTemples			
			Time ↓	Size(MB) ↓	PSNR ↑	SSIM ↑	PSNR ↑	SSIM ↑	PSNR ↑	SSIM ↑
SRN [36]	-	-	>10h	-	22.26	0.846	24.33	0.882	24.10	0.847
NSVF [18]	8192	150k	>48*h	-	31.75	0.953	35.18	0.979	28.48	0.901
NeRF [24]	4096	300k	~35h	5.00	31.01	0.947	30.81	0.952	25.78	0.864
SNeRG [12]	8192	250k	~15h	1771.5	30.38	0.950	-	-	-	-
PlenOctrees [47]	1024	200k	~15h	1976.3	31.71	0.958	-	-	27.99	0.917
Plenoxels [46]	5000	128k	11.4m	778.1	31.71	0.958	-	-	27.43	0.906
DVGO [37]	5000	30k	15.0m	612.1	31.95	0.957	35.08	0.975	28.41	0.911
Ours-CP-384	4096	30k	25.2m	3.9	31.56	0.949	34.48	0.971	27.59	0.897
Ours-VM-192-SH	4096	30k	16.8m	71.9	32.00	0.955	35.30	0.977	27.81	0.907
Ours-VM-48	4096	30k	13.8m	18.9	32.39	0.957	35.34	0.976	28.06	0.909
Ours-VM-192	4096	15k	8.1m	71.8	32.52	0.959	35.59	0.978	28.07	0.913
Ours-VM-192	4096	30k	17.4m	71.8	33.14	0.963	36.52	0.982	28.56	0.920

Table 1: We compare our method with previous and concurrent novel view synthesis methods on three datasets. All scores of the baseline methods are directly taken from their papers whenever available. We also report the averaged reconstruction time and model size for the Synthetic-NeRF dataset. NSVF requires 8 GPUs for optimization (marked by a star), while others run on a single GPU. DVGO’s 30k steps correspond to 10k for coarse and 20k for fine reconstruction.

is able to outperform the best CP model that uses many more components and voxels. Remarkably, the PSNR of 31.81 from this smallest VM model (which only takes 8.6 MB) is already higher than the PSNRs of NeRF and many other previous and concurrent techniques (see Tab. 1). In addition, 192 components are generally adequate for TensoRF-VM; doubling the number to 384 only leads to marginal improvement. TensoRF-VM with 300^3 voxels can already lead to high PSNRs close to or greater than 33, while retaining compact model sizes (<72MB). Increasing the resolution further leads to improved quality, but also larger model size.

Also note that all of our TensoRF models can achieve very fast reconstruction. Except for the largest VM model, all models finish reconstruction in less than 30 min, significantly faster than NeRF and many previous methods (see Tab. 1).

Optimization steps. We further evaluate our approach with different optimization steps for our best CP model and the VM models with 300^3 voxels. PSNRs and reconstruction time are shown in Tab. 3. All of our models consistently achieve better rendering quality with more steps. Our compact CP-384 model (3.9MB) can even achieve a PSNR greater than 32 after 60k steps, higher than the PSNRs of all previous methods in Tab. 1. On the other hand, our VM models can quickly achieve high rendering quality in very few steps. With only 15k steps, many models achieve high PSNRs that are already state-of-the-art.

Comparisons on 360° scenes. We compare our approach with state-of-the-art novel view synthesis methods, including previous works (SRN[36], NeRF[24], NSVF[18]), SNeRG[12], PlenOctrees[47]) and concurrent works (Plenoxels [46], DVGO[37]). In particular, we compare with them using our best CP model and our VM models (300^3 voxels) with 48 and 192 components. We also show a 192-component VM model with spherical harmonics shading function. Table

	#Comp	200 ³	300 ³	500 ³
TensorRF-CP	48	27.98/09:29/0.74	28.24/11:45/1.09	28.38/14:20/1.85
	96	28.50/09:57/0.88	28.83/12:12/1.29	29.06/15:27/2.18
	192	29.50/11:09/1.08	29.99/13:41/1.59	30.33/18.03/2.66
	384	30.47/14:41/1.59	31.08/18:09/2.33	31.56/25:11/3.93
TensorRF-VM	48	31.81/11:29/08.6	32.39/13:51/23.5	32.63/18:17/55.8
	96	32.33/11:54/16.5	32.86/14:08/37.3	33.06/20:11/105.
	192	32.63/13:26/32.3	33.14/17:36/76.7	33.31/27.18/204.
	384	32.69/17:24/63.4	33.21/25:14/143.	33.39/43:19/397.

Table 2: We compare the averaged PSNRs / optimization time (mm:ss) / model sizes (MB) of CP and VM TensorRF models on Synthetic NeRF dataset [24] with different numbers of components and grid resolutions, optimized for 30k steps.

	5k	15k	30k	60k	5k	15k	30k	60k
CP-384	28.37	30.80	31.56	32.03	03:03	11:30	25:11	51:47
VM-48	29.28	31.80	32.39	32.68	01:57	06:21	13:51	27:20
VM-96	29.65	32.26	32.86	33.17	02:01	06:41	14:08	28:57
VM-192	29.86	32.52	33.14	33.44	02:16	08:08	17:37	35:50
VM-384	29.95	32.62	33.21	33.52	02:51	11:30	25:14	52:50

Table 3: PSNRs and time of CP and VM models with different training steps on the Synthetic-NeRF dataset [24].

Method	Time ↓	Size	PSNR↑	SSIM↑
NeRF [24]	36h	5.00M	26.50	0.811
Plenoxels [46]	24:00m	2.59G	26.29	0.829
Ours-VM-48	19:44m	90.4M	26.51	0.832
Ours-VM-96	25.43m	179.7M	26.73	0.839

Table 4: Quantitative comparisons of our method with NeRF and Plenoxels on forward-facing scenes [23].

1 shows the quantitative results (PSNRs and SSIMs) of ours and comparison methods on three challenging datasets, where we also show the corresponding batch size, optimization steps, time, and final output model size for each model, to compare all methods from multiple perspectives. Note that all of our CP and VM models can outperform NeRF on all three datasets while taking substantially less optimization time and fewer steps. Our best VM-192 model can even achieve the best PSNRs and SSIMs on all datasets, significantly outperforming the comparison methods. Our approach can also achieve qualitatively better renderings with more appearance and geometry details and less outliers, as shown in Fig. 4.

Our models are highly efficient, which all require less than 75MB space and can be reconstructed in less than 30 min. This corresponds to more than 70x speed up compared to NeRF that requires about 1.5 days for optimization. Our CP model is even more compact than NeRF. Moreover, SNeRG and PlenOctrees require pre-training a NeRF-like MLP, requiring long reconstruction time too. DVGO and Plenoxels are concurrent works, which can also achieve fast reconstruction in less than 15 min. However, as both are voxel-based methods and directly optimize voxel values, they lead to huge model sizes similar to previous voxel-based methods like SNeRG and PlenOctrees. In contrast, we factorize feature grids and model them with compact vectors and matrices, leading to substantially smaller model sizes. Meanwhile, our VM-192 can even reconstruct faster than DVGO and Plenoxels, taking only 15k steps, and achieving better quality in most cases. In fact, Plenoxels’ fast reconstruction relies on quickly optimizing significantly more steps (> 4 times our steps) with their CUDA implementation. Our models are implemented with standard PyTorch modules and already achieve much better rendering quality with fewer steps taking comparable and even less reconstruction time than Plenoxels. Note that our SH model essentially represents

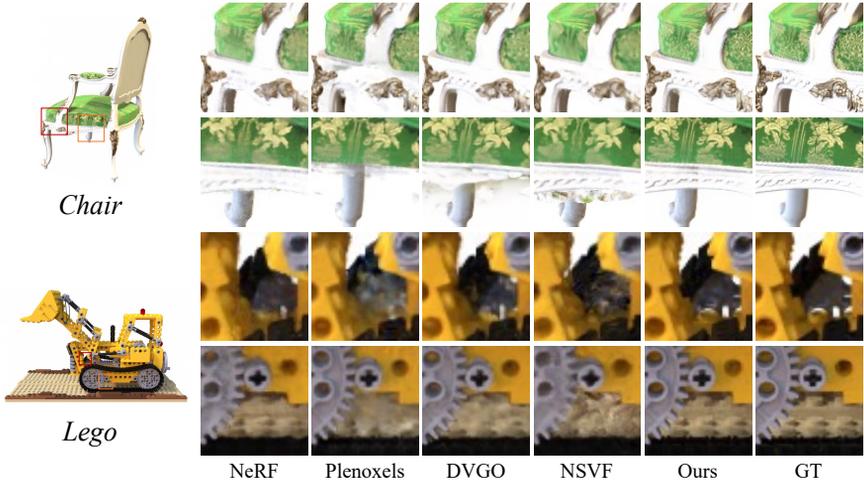


Fig. 4: Qualitative results of our VM-192-30k model and comparison methods (NeRF[24], plenoxels [46], DVG0 [37], NSVF [18]) on two Synthetic NeRF scenes.

the same underlying feature grid as Plenoxels but can still lead to more compact modeling and better quality with fewer steps, showing the advantages of our factorization based modeling. In general, our approach enables fast reconstruction, compact modeling, and photo-realistic rendering simultaneously.

Forward-facing scenes. Our approach can also achieve efficient and high-quality radiance field reconstruction for forward-facing scenes. We show quantitative results of our two VM models on the LLFF dataset [23] and compare with NeRF and Plenoxels in Tab. 4. Our models outperform the previous state-of-the-art NeRF and take significantly less reconstruction time. Compared with Plenoxels [46], our approach leads to comparable or faster reconstruction speed, better quality, and substantially smaller model sizes.

7 Conclusion.

We have presented a novel approach for high-quality scene reconstruction and rendering. We propose a novel scene representation – TensoRF which leverages tensor decomposition techniques to model radiance fields compactly as factorized low-rank tensor components. While our framework accommodates classical tensor factorization techniques (like CP), we introduce a novel vector-matrix decomposition that leads to better reconstruction quality and faster optimization speed. Our approach enables highly efficient radiance field reconstruction in less than 30 min per scene, leading to better rendering quality compared to NeRF that requires substantially longer training time (20+ hours). Moreover, our tensor factorization based method achieves high compactness, leading to a memory footprint of less than 75MB, substantially smaller than many other previous and concurrent voxel-grid based methods.

References

1. Ben-Younes, H., Cadene, R., Thome, N., Cord, M.: Block: Bilinear superdiagonal fusion for visual question answering and visual relationship detection. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 33, pp. 8102–8109 (2019)
2. Bi, S., Xu, Z., Srinivasan, P., Mildenhall, B., Sunkavalli, K., Hašan, M., Hold-Geoffroy, Y., Kriegman, D., Ramamoorthi, R.: Neural reflectance fields for appearance acquisition. arXiv preprint arXiv:2008.03824 (2020)
3. Bi, S., Xu, Z., Sunkavalli, K., Hašan, M., Hold-Geoffroy, Y., Kriegman, D., Ramamoorthi, R.: Deep reflectance volumes: Relightable reconstructions from multi-view photometric images. In: Proc. ECCV (2020)
4. Boss, M., Braun, R., Jampani, V., Barron, J.T., Liu, C., Lensch, H.: Nerd: Neural reflectance decomposition from image collections. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 12684–12694 (2021)
5. Carroll, J.D., Chang, J.J.: Analysis of individual differences in multidimensional scaling via an n-way generalization of “eckart-young” decomposition. *Psychometrika* **35**(3), 283–319 (1970)
6. Chan, E.R., Monteiro, M., Kellnhofer, P., Wu, J., Wetzstein, G.: pi-gan: Periodic implicit generative adversarial networks for 3d-aware image synthesis. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5799–5809 (2021)
7. Chen, A., Xu, Z., Zhao, F., Zhang, X., Xiang, F., Yu, J., Su, H.: Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. arXiv preprint arXiv:2103.15595 (2021)
8. Chibane, J., Bansal, A., Lazova, V., Pons-Moll, G.: Stereo radiance fields (srf): Learning view synthesis from sparse views of novel scenes. In: IEEE Conference on Computer Vision and Pattern Recognition (CVPR). IEEE (jun 2021)
9. De Lathauwer, L.: Decompositions of a higher-order tensor in block terms—part ii: Definitions and uniqueness. *SIAM Journal on Matrix Analysis and Applications* **30**(3), 1033–1066 (2008)
10. Groueix, T., Fisher, M., Kim, V.G., Russell, B.C., Aubry, M.: A papier-mâché approach to learning 3d surface generation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 216–224 (2018)
11. Harshman, R.A., et al.: Foundations of the parafac procedure: Models and conditions for an “ explanatory” multimodal factor analysis (1970)
12. Hedman, P., Srinivasan, P.P., Mildenhall, B., Barron, J.T., Debevec, P.: Baking neural radiance fields for real-time view synthesis. arXiv preprint arXiv:2103.14645 (2021)
13. Ji, M., Gall, J., Zheng, H., Liu, Y., Fang, L.: SurfaceNet: An end-to-end 3D neural network for multiview stereopsis. In: Proc. ICCV (2017)
14. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014)
15. Knapitsch, A., Park, J., Zhou, Q.Y., Koltun, V.: Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics* **36**(4) (2017)
16. Kolda, T.G., Bader, B.W.: Tensor decompositions and applications. *SIAM review* **51**(3), 455–500 (2009)
17. Li, Z., Niklaus, S., Snavely, N., Wang, O.: Neural scene flow fields for space-time view synthesis of dynamic scenes. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6498–6508 (2021)

18. Liu, L., Gu, J., Lin, K.Z., Chua, T.S., Theobalt, C.: Neural sparse voxel fields. arXiv preprint arXiv:2007.11571 (2020)
19. Liu, S., Zhang, X., Zhang, Z., Zhang, R., Zhu, J.Y., Russell, B.: Editing conditional radiance fields. arXiv preprint arXiv:2105.06466 (2021)
20. Lombardi, S., Simon, T., Saragih, J., Schwartz, G., Lehrmann, A., Sheikh, Y.: Neural volumes: Learning dynamic renderable volumes from images. arXiv preprint arXiv:1906.07751 (2019)
21. Martin-Brualla, R., Radwan, N., Sajjadi, M.S., Barron, J.T., Dosovitskiy, A., Duckworth, D.: Nerf in the wild: Neural radiance fields for unconstrained photo collections. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7210–7219 (2021)
22. Mescheder, L., Oechsle, M., Niemeyer, M., Nowozin, S., Geiger, A.: Occupancy networks: Learning 3d reconstruction in function space. Proc. CVPR (2019)
23. Mildenhall, B., Srinivasan, P.P., Ortiz-Cayon, R., Kalantari, N.K., Ramamoorthi, R., Ng, R., Kar, A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. ACM Transactions on Graphics (TOG) **38**(4), 1–14 (2019)
24. Mildenhall, B., Srinivasan, P.P., Tancik, M., Barron, J.T., Ramamoorthi, R., Ng, R.: Nerf: Representing scenes as neural radiance fields for view synthesis. In: European conference on computer vision. pp. 405–421. Springer (2020)
25. Müller, T., Evans, A., Schied, C., Keller, A.: Instant neural graphics primitives with a multiresolution hash encoding. arXiv preprint arXiv:2201.05989 (2022)
26. Niemeyer, M., Geiger, A.: Giraffe: Representing scenes as compositional generative neural feature fields. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 11453–11464 (2021)
27. Oechsle, M., Peng, S., Geiger, A.: Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction. arXiv preprint arXiv:2104.10078 (2021)
28. Park, K., Sinha, U., Hedman, P., Barron, J.T., Bouaziz, S., Goldman, D.B., Martin-Brualla, R., Seitz, S.M.: Hypernerf: A higher-dimensional representation for topologically varying neural radiance fields. arXiv preprint arXiv:2106.13228 (2021)
29. Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z., Gimelshein, N., Antiga, L., et al.: Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems **32** (2019)
30. Peng, S., Niemeyer, M., Mescheder, L., Pollefeys, M., Geiger, A.: Convolutional occupancy networks. In: European Conference on Computer Vision. pp. 523–540. Springer (2020)
31. Qi, C.R., Su, H., Mo, K., Guibas, L.J.: Pointnet: Deep learning on point sets for 3d classification and segmentation. In: Proc. CVPR (2017)
32. Qi, C.R., Su, H., Nießner, M., Dai, A., Yan, M., Guibas, L.J.: Volumetric and multi-view cnns for object classification on 3d data. In: Proc. CVPR (2016)
33. Reiser, C., Peng, S., Liao, Y., Geiger, A.: Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. arXiv preprint arXiv:2103.13744 (2021)
34. Schwarz, K., Liao, Y., Niemeyer, M., Geiger, A.: Graf: Generative radiance fields for 3d-aware image synthesis. arXiv preprint arXiv:2007.02442 (2020)
35. Sitzmann, V., Thies, J., Heide, F., Nießner, M., Wetzstein, G., Zollhofer, M.: Deepvoxels: Learning persistent 3D feature embeddings. In: Proc. CVPR (2019)
36. Sitzmann, V., Zollhöfer, M., Wetzstein, G.: Scene representation networks: Continuous 3d-structure-aware neural scene representations. arXiv preprint arXiv:1906.01618 (2019)
37. Sun, C., Sun, M., Chen, H.T.: Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. arXiv preprint arXiv:2111.11215 (2021)

38. Tucker, L.R.: Some mathematical notes on three-mode factor analysis. *Psychometrika* **31**(3), 279–311 (1966)
39. Wang, N., Zhang, Y., Li, Z., Fu, Y., Liu, W., Jiang, Y.G.: Pixel2mesh: Generating 3d mesh models from single RGB images. In: Proc. ECCV (2018)
40. Wang, P., Liu, L., Liu, Y., Theobalt, C., Komura, T., Wang, W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. arXiv preprint arXiv:2106.10689 (2021)
41. Wang, Q., Wang, Z., Genova, K., Srinivasan, P., Zhou, H., Barron, J.T., Martin-Brualla, R., Snavely, N., Funkhouser, T.: Ibrnet: Learning multi-view image-based rendering. In: CVPR (2021)
42. Xiang, F., Xu, Z., Hasan, M., Hold-Geoffroy, Y., Sunkavalli, K., Su, H.: Neutex: Neural texture mapping for volumetric neural rendering. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 7119–7128 (2021)
43. Xu, Q., Xu, Z., Philip, J., Bi, S., Shu, Z., Sunkavalli, K., Neumann, U.: Point-nerf: Point-based neural radiance fields. arXiv preprint arXiv:2201.08845 (2022)
44. Ye, J., Li, G., Chen, D., Yang, H., Zhe, S., Xu, Z.: Block-term tensor neural networks. *Neural Networks* **130**, 11–21 (2020)
45. Ye, J., Wang, L., Li, G., Chen, D., Zhe, S., Chu, X., Xu, Z.: Learning compact recurrent neural networks with block-term tensor decomposition. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 9378–9387 (2018)
46. Yu, A., Fridovich-Keil, S., Tancik, M., Chen, Q., Recht, B., Kanazawa, A.: Plenoxels: Radiance fields without neural networks. arXiv preprint arXiv:2112.05131 (2021)
47. Yu, A., Li, R., Tancik, M., Li, H., Ng, R., Kanazawa, A.: Plenotrees for real-time rendering of neural radiance fields. arXiv preprint arXiv:2103.14024 (2021)
48. Yu, A., Ye, V., Tancik, M., Kanazawa, A.: pixelnerf: Neural radiance fields from one or few images. In: CVPR (2021)
49. Zhang, K., Riegler, G., Snavely, N., Koltun, V.: Nerf++: Analyzing and improving neural radiance fields. arXiv preprint arXiv:2010.07492 (2020)
50. Zhou, T., Tucker, R., Flynn, J., Fyffe, G., Snavely, N.: Stereo magnification: learning view synthesis using multiplane images. *ACM Transactions on Graphics* **37**(4), 1–12 (2018)

A TensorRF Representation Details.

We illustrate the feature grid of our tensorial radiance field and the tensor factors in TensorRF with both CP and VM decompositions in Fig. 5.

Number of components. The total number of tensor components ($\#Comp$) is $(R_\sigma + R_c)$ for TensorRF-CP and $3(R_\sigma + R_c)$ for TensorRF-VM (because VM has three types of components). Therefore, the R we use for TensorRF-CP is three times as large as the R used for TensorRF-VM to achieve the same number of components shown in Tab. 2. We also find that using $R_\sigma < R_c$ is usually better than $R_\sigma = R_c$ when R_σ is large enough (> 8). In particular, for TensorRF-VM, we use $R_\sigma = R_c = 8$ for $\#Comp = 48$; $R_\sigma = 8, R_c = 24$ for $\#Comp = 96$; $R_\sigma = 16, R_c = 48$ for $\#Comp = 192$; $R_\sigma = 32, R_c = 96$ for $\#Comp = 384$. Note that, as discussed in Eqn. 3,4, here we apply the same number of components for $\mathcal{A}^X, \mathcal{A}^Y, \mathcal{A}^Z$ with $R_1 = R_2 = R_3 = R$ for both density and appearance (where R is R_σ and R_c respectively), assuming the three spatial dimensions are equally complex.

Forward-facing settings. We use the above settings with $R_1 = R_2 = R_3$, for all 360° object datasets in Tab. 1. On the other hand, Forward-facing scenes apparently appear differently in the three dimensions; especially, in NDC space, the X and Y spatial modes (corresponding to the image plane) contains more appearance information that is visible to rendering viewpoints. We therefore use more components for the $X - Y$ plane, corresponding to $\mathcal{A}^Z = \mathbf{v}^Z \circ \mathbf{M}^{X,Y}$. In this case, these \mathcal{A}^Z components can also be seen as special compressed versions of neural MPIs. In particular, the detailed numbers of components we use for generating the results in Tab. 4 are: for $\#Comp=48$ $R_{\sigma,1} = R_{\sigma,2} = 4, R_{\sigma,3} = 16, R_{c,1} = R_{c,2} = 4, R_{c,3} = 16$; for $\#Comp=96$, $R_{\sigma,1} = R_{\sigma,2} = 4, R_{\sigma,3} = 16, R_{c,1} = R_{c,2} = 16, R_{c,3} = 16$.

Number of parameters. We briefly discuss the number of parameters in our model. With the same $\#Comp$, when $I = J = K$ and $R_\sigma + R_c = R$, the total number of parameters used for TensorRF-CP is $3KR + PR_c$; for Tensor-VM, the number is $K^2R + KR + PR_c$ (here considering $R_\sigma/3, R_c/3$ are used to make the $\#Comp$ the same as TensorRF-CP). For example, for a $300 \times 300 \times 300$ feature grid with $P = 27$ channels (plus one density channel), the total number of parameters in a dense grid is 756 M; the number of parameters used for TensorRF-CP (when $R = 192$) is 360 K; the number of parameters used for TensorRF-VM (when $R = 192$) is 17 M. Our CP and VM model can achieve 0.048% and 2.25% compression rates respectively.

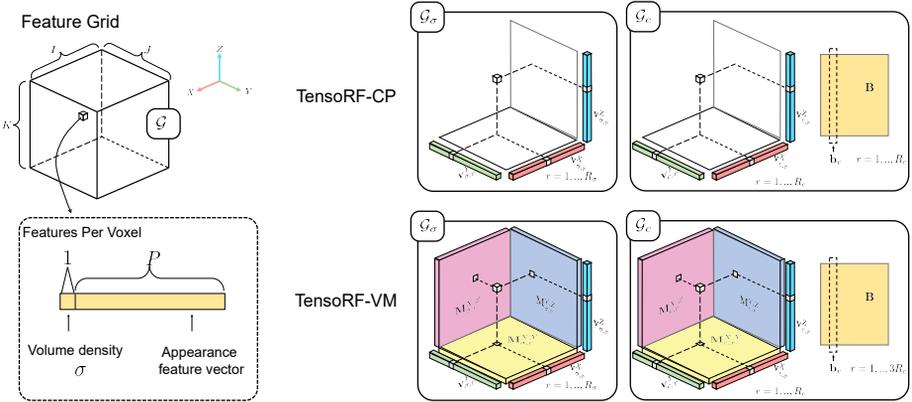


Fig. 5: Feature grids and factorized tensors in TensorRF. We leverage a regular voxel grid \mathcal{G} , covering a 3D scene, to model a radiance field of the scene. Each voxel of \mathcal{G} contains multi-channel features, where one channel represents the volume density (σ) and the remaining P channels lead to an appearance feature vector (f_c) for computing view-dependent colors. We split the density and appearance features into two feature grids \mathcal{G}_σ and \mathcal{G}_c , consider them as 3D and 4D tensors, and factorize them into compact factors with outer products. TensorRF with CP decomposition factorize the tensors into only vectors and TensorRF with VM decomposition factorize the tensor into vector and matrix factors (Eqn. 7, 8). Note that, each voxel of the original grid is only related to one value from each XYZ-mode vector/matrix factor in both decompositions, which are marked in the figure.

B More Implementation Details.

Loss functions. As described in sec. 4.4, we apply a L2 rendering loss and additional regularization terms to optimize our tensor factors for radiance field reconstruction. In general, this loss function is expressed by

$$\mathcal{L} = \|C - \tilde{C}\|_2^2 + \omega \cdot \mathcal{L}_{reg} \quad (17)$$

Where \tilde{C} is the groundtruth color, ω is weight of the regularization term.

To encourage the sparsity in the parameters of our tensor factors, we apply the standard L1 regularization, which we find is effective in improving the quality in extrapolating views and removing floaters/outliers in final renderings. Note that, unlike previous methods [12,37] that penalize predicted per-point density with a Cauchy loss or entropy loss, our L1 regularizer is much simpler and directly applied on the parameters of tensor factors. We find that it is sufficient to apply the L1 sparsity loss only on the density parameters, expressed by

$$\mathcal{L}_{L1} = \frac{1}{N} \sum_{r=1}^{R_\sigma} (\|\mathbf{M}_{\sigma,r}\| + \|\mathbf{v}_{\sigma,r}\|), \quad (18)$$

where $\|\mathbf{M}_{\sigma,r}\|$ and $\|\mathbf{v}_{\sigma,r}\|$ are simply the sum of absolute values of all elements, and $N = R_\sigma \cdot (I \cdot J + I \cdot K + J \cdot K + I + J + K)$ is the total number of parameters. We use this L1 sparsity loss with a $\omega = 0.0004$ for the Synthetic NeRF and Synthetic NSVF datasets. An ablation study on this L1 loss on the Synthetic NeRF dataset is shown in Tab. .

For real datasets that have very few input images (like LLFF[23]) or imperfect capture conditions (like Tanks and Temples [15] that has varying exposure and inconsistent masks), we find a TV loss is more efficient than the L1 sparsity loss, expressed by

$$\mathcal{L}_{TV} = \frac{1}{N} \sum (\sqrt{\Delta^2 \mathcal{A}_{\sigma,r}^m} + 0.1 \cdot \sqrt{\Delta^2 \mathcal{A}_{C,r}^m}), \quad (19)$$

Here Δ^2 is the squared difference between the neighboring values in the matrix/vector factors; we apply a smaller weight (weighted by 0.1 additionally) on appearance parameters in the TV loss. We use $\omega = 1$ when using this TV loss. **More details.** As described in Sec. 5, we use a small two-layer MLP with 128 channels in hidden layers as our neural decoding function. In particular, the input to this MLP contains the viewing direction and the features recovered by our tensor factors (no xyz positions are used). Similar to NeRF and NSVF [24,18], we also apply frequency encodings (with Sin and Cos functions) on both the viewing direction and features. Unlike NeRF that uses ten different frequencies, we use only two.

During optimization, we also apply an exponential learning rate decay to make the optimization more stable when the reconstruction is being finished. Specifically, we decay our initial learning rates at every training step, until decayed by a factor of 0.1 in the end of the optimization.

C More Evaluation.

We perform an ablation study to evaluate our L1 regularization. Tab. 5 shows how our framework performs by removing the L1 regularization on the Synthetic-NeRF dataset, our models exceed NeRF fidelity (31.01 in average) even without regularization. We observe the performance gap between w/ and w/o L1 regularization is mostly caused by the floaters in the empty space. We also provide more results on our model with different numbers of training steps in Tab. 6, which is basically a detailed version of Tab. 3 with more settings. These results showcase that our models consistently improve when training with more iterations.

D Limitations and Future Work.

Our approach achieves high-quality radiance field reconstruction for 360° objects and forward-facing scenes; however, our method currently only supports bounded scenes with a single bounding box and cannot handle unbounded scenes with both foreground and background content. Combining our method with techniques

	PSNR	SSIM	LPIPS _{VGG}	LPIPS _{Alex}
CP-384	31.56/31.23	0.949/0.947	0.076/0.078	0.041/0.043
VM-48	32.39/31.71	0.957/0.953	0.057/0.062	0.032/0.036
VM-192-SH	32.00/31.14	0.955/0.949	0.058/0.068	0.058/0.044
VM-192	33.14/32.43	0.963/0.960	0.047/0.052	0.027/0.030

Table 5: We compare the averaged scores against w/o L1 regularization on the Synthetic-NeRF dataset.

	5k	8k	10k	12k	15k	20k	30k	40k	60k	100k
CP-192	28.38	29.13	29.93	30.38	30.80	31.18	31.56	31.75	32.03	32.18
VM-48	29.28	30.39	31.11	31.47	31.80	32.08	32.39	32.55	32.68	32.84
VM-96	29.65	30.72	31.52	31.93	32.26	32.56	32.86	33.00	33.17	33.29
VM-192	29.86	30.93	31.74	32.17	32.52	32.85	33.14	33.27	33.44	33.54
VM-384	29.95	30.88	31.75	32.20	32.62	32.94	33.21	33.35	33.52	33.64

Table 6: PSNRs on the Synthetic NeRF datasets with different numbers of training steps. This is more detailed version than Tab. 3.

like NeRF++ [49] to separately model a foreground field inside a regular box and a background field inside another box defined in a spherical coordinate space can potentially extend our method to address unbounded scenes. Despite the success in per-scene optimization shown in this paper, an interesting future direction is to discover and learn general basis factors across scenes on a large-scale dataset, leveraging data priors to further improve the quality or enable other applications like GANs (as done in GRAF [34] and GIRRAF [26]).

E Acknowledgements

We would like to thank Yannick Hold-Geoffroy for his useful tips in video animation, Qiangeng Xu for providing some baseline results, and, Katja Schwarz and Michael Niemeyer for providing helpful video materials. This project was supported by NSF grant IIS-1764078 and gift money from VIVO.

F Per-scene Breakdown.

Tab. 8-11 provide a per-scene break down for quantity metrics in Synthesis-nerf [21], Synthe-nsvf [18], Tanks&Templates [15] and forward-facing [23] dataset.

Methods	Avg.	<i>Chair</i>	<i>Drums</i>	<i>Ficus</i>	<i>Hotdog</i>	<i>Lego</i>	<i>Materials</i>	<i>Mic</i>	<i>Ship</i>
PSNR\uparrow									
SRN [36]	22.26	26.96	17.18	20.73	26.81	20.85	18.09	26.85	20.60
NeRF [24]	31.01	33.00	25.01	30.13	36.18	32.54	29.62	32.91	28.65
NSVF [18]	31.75	33.19	25.18	31.23	37.14	32.29	32.68	34.27	27.93
SNeRG [12]	30.38	33.24	24.57	29.32	34.33	33.82	27.21	32.60	27.97
PlenOctrees [47]	31.71	34.66	25.31	30.79	36.79	32.95	29.76	33.97	29.42
Plenoxels [46]	31.71	33.98	25.35	31.83	36.43	34.10	29.14	33.26	29.62
DVGO [37]	31.95	34.09	25.44	32.78	36.74	34.64	29.57	33.20	29.13
Ours-CP-384	31.56	33.60	25.17	30.72	36.24	34.05	30.10	33.77	28.84
Ours-VM-192-SH	32.00	34.68	25.37	32.30	36.30	35.42	29.30	33.21	29.46
Ours-VM-48	32.39	34.68	25.58	33.37	36.81	35.51	29.45	33.59	30.12
Ours-VM-192-15k	32.52	34.95	25.63	33.46	36.85	35.78	29.78	33.69	30.04
Ours-VM-192-30k	33.14	35.76	26.01	33.99	37.41	36.46	30.12	34.61	30.77

Table 7: PSNR results on each scene from the **Synthetic-NeRF** [24] dataset.

Methods	Avg.	<i>Chair</i>	<i>Drums</i>	<i>Ficus</i>	<i>Hotdog</i>	<i>Lego</i>	<i>Materials</i>	<i>Mic</i>	<i>Ship</i>
SSIM ↑									
SRN [36]	0.846	0.910	0.766	0.849	0.923	0.809	0.808	0.947	0.757
NeRF [24]	0.947	0.967	0.925	0.964	0.974	0.961	0.949	0.980	0.856
NSVF [18]	0.953	0.968	0.931	0.973	0.980	0.960	0.973	0.987	0.854
SNeRG [12]	0.950	0.975	0.929	0.967	0.971	0.973	0.938	0.982	0.865
PlenOctrees [47]	0.958	0.981	0.933	0.970	0.982	0.971	0.955	0.987	0.884
Plenoxels [46]	0.958	0.977	0.933	0.976	0.980	0.976	0.949	0.985	0.890
DVGO [37]	0.957	0.977	0.930	0.978	0.980	0.976	0.951	0.983	0.879
Ours-CP-384	0.949	0.973	0.921	0.965	0.975	0.971	0.950	0.983	0.857
Ours-VM-192-SH	0.955	0.979	0.928	0.976	0.977	0.978	0.941	0.983	0.875
Ours-VM-48	0.957	0.980	0.929	0.979	0.979	0.979	0.942	0.984	0.883
Ours-VM-192-15k	0.959	0.982	0.933	0.981	0.980	0.981	0.949	0.985	0.886
Ours-VM-192-30k	0.963	0.985	0.937	0.982	0.982	0.983	0.952	0.988	0.895
LPIPS _{VGG} ↓									
SRN [36]	0.170	0.106	0.267	0.149	0.100	0.200	0.174	0.063	0.299
NeRF [24]	0.081	0.046	0.091	0.044	0.121	0.050	0.063	0.028	0.206
PlenOctrees [47]	0.053	0.022	0.076	0.038	0.032	0.034	0.059	0.017	0.144
Plenoxels [46]	0.049	0.031	0.067	0.026	0.037	0.028	0.057	0.015	0.134
DVGO [37]	0.053	0.027	0.077	0.024	0.034	0.028	0.058	0.017	0.161
Ours-CP-384	0.076	0.044	0.114	0.058	0.052	0.038	0.068	0.035	0.196
Ours-VM-192-SH	0.058	0.031	0.082	0.028	0.048	0.024	0.069	0.022	0.160
Ours-VM-48	0.057	0.030	0.087	0.028	0.039	0.024	0.072	0.021	0.155
Ours-VM-192-15k	0.053	0.026	0.078	0.025	0.038	0.021	0.063	0.020	0.153
Ours-VM-192-30k	0.047	0.022	0.073	0.022	0.032	0.018	0.058	0.015	0.138
LPIPS _{Alex} ↓									
NSVF [18]	0.047	0.043	0.069	0.017	0.025	0.029	0.021	0.010	0.162
DVGO [37]	0.035	0.016	0.061	0.015	0.017	0.014	0.026	0.014	0.118
Ours-CP-384	0.041	0.022	0.069	0.024	0.024	0.014	0.031	0.018	0.130
Ours-VM-192-SH	0.058	0.031	0.082	0.028	0.048	0.024	0.069	0.022	0.160
Ours-VM-48	0.032	0.014	0.059	0.015	0.017	0.009	0.036	0.012	0.098
Ours-VM-192-15k	0.032	0.013	0.056	0.014	0.017	0.009	0.029	0.013	0.101
Ours-VM-192-30k	0.027	0.010	0.051	0.012	0.013	0.007	0.026	0.009	0.085

Table 8: Quantitative results on each scene from the **Synthetic-NeRF** [24] dataset.

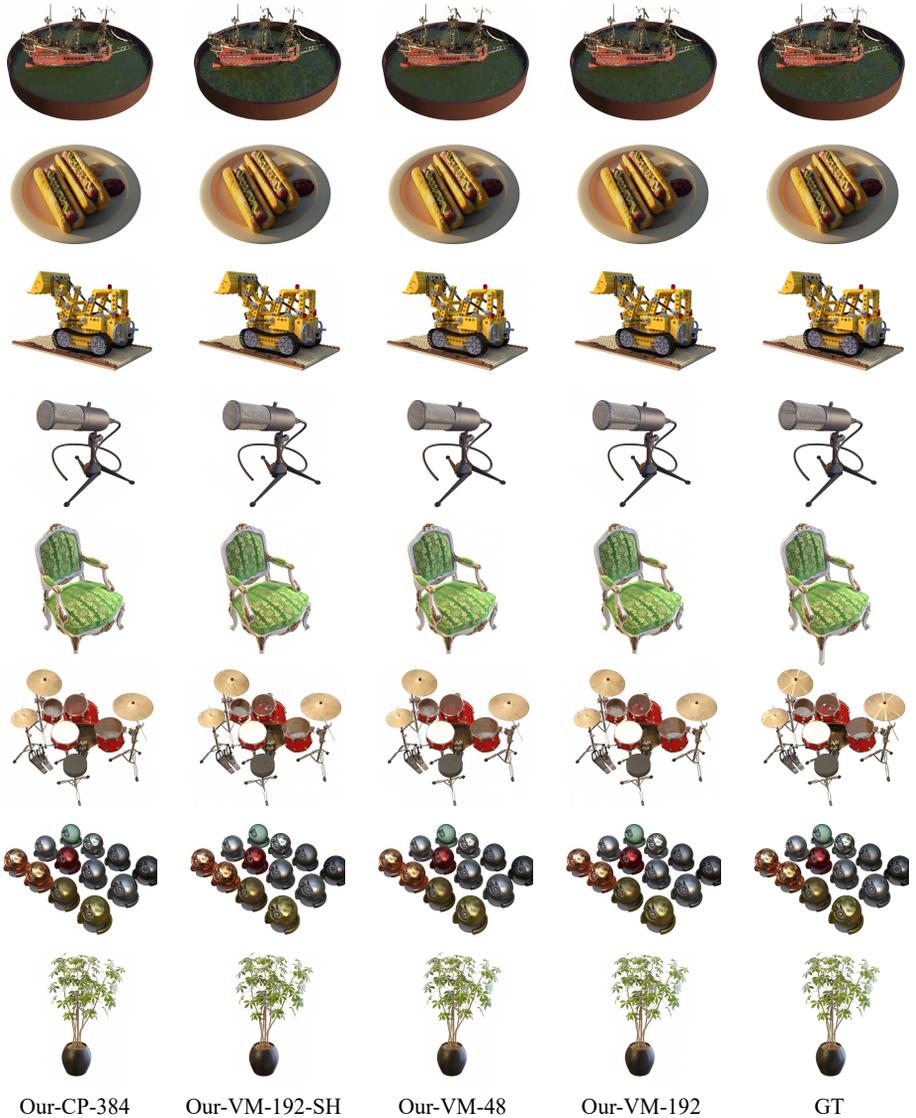


Fig. 6: Our rendering results on **Synthetic-NeRF** dataset. From top to bottom: Ship, Hotdog, Lego, Mic, Chair, Drums, Materials, Ficus.

Methods	Avg.	Winholder	Steamtrain	Toad	Robot	Bike	Palace	Spaceship	Lifestyle
PSNR\uparrow									
SRN [36]	24.33	20.74	25.49	25.36	22.27	23.76	24.45	27.99	24.58
NeRF [24]	30.81	28.23	30.84	29.42	28.69	31.77	31.76	34.66	31.08
NSVF [18]	35.13	32.04	35.13	33.25	35.24	37.75	34.05	39.00	34.60
DVGO [37]	35.08	30.26	36.56	33.10	36.36	38.33	34.49	37.71	33.79
Ours-CP-384	34.48	29.92	36.07	31.37	35.92	36.74	36.26	37.01	32.54
Ours-VM-192-SH	35.30	29.72	37.33	34.03	37.59	38.61	36.09	35.82	33.21
Ours-VM-48	35.34	30.46	37.06	33.13	36.92	37.98	36.32	37.19	33.68
Ours-VM-192-15k	35.59	30.31	37.20	33.63	37.29	38.33	36.57	37.77	33.62
Ours-VM-192-30k	36.52	31.32	37.87	34.85	38.26	39.23	37.56	38.60	34.51
SSIM\uparrow									
SRN [36]	0.882	0.850	0.923	0.822	0.904	0.926	0.792	0.945	0.892
NeRF [24]	0.952	0.920	0.966	0.920	0.960	0.970	0.950	0.980	0.946
NSVF [18]	0.979	0.965	0.986	0.968	0.988	0.991	0.969	0.991	0.971
DVGO [37]	0.975	0.949	0.989	0.966	0.992	0.991	0.962	0.988	0.965
Ours-CP-384	0.971	0.947	0.986	0.950	0.990	0.987	0.971	0.984	0.951
Ours-VM-192-SH	0.977	0.953	0.988	0.974	0.993	0.991	0.972	0.982	0.964
Ours-VM-48	0.976	0.952	0.988	0.968	0.992	0.990	0.973	0.985	0.962
Ours-VM-192-15k	0.978	0.953	0.989	0.972	0.993	0.991	0.975	0.987	0.964
Ours-VM-192-30k	0.982	0.961	0.991	0.978	0.994	0.993	0.979	0.989	0.968
LPIPS_{VGG} \downarrow									
DVGO [37]	0.033	0.055	0.019	0.047	0.013	0.011	0.043	0.019	0.054
Ours-CP-384	0.045	0.082	0.031	0.067	0.016	0.023	0.031	0.028	0.084
Ours-VM-192-SH	0.031	0.057	0.024	0.035	0.011	0.013	0.030	0.026	0.051
Ours-VM-48	0.034	0.061	0.023	0.047	0.013	0.014	0.029	0.025	0.059
Ours-VM-192-15k	0.031	0.060	0.020	0.040	0.011	0.012	0.028	0.022	0.055
Ours-VM-192-30k	0.026	0.051	0.017	0.031	0.010	0.010	0.022	0.020	0.048
LPIPS_{Alex} \downarrow									
SRN [36]	0.141	0.224	0.082	0.204	0.120	0.075	0.240	0.061	0.120
NeRF [24]	0.043	0.096	0.031	0.069	0.038	0.019	0.031	0.016	0.047
NSVF [18]	0.015	0.020	0.010	0.032	0.007	0.004	0.018	0.006	0.020
DVGO [37]	0.019	0.038	0.010	0.030	0.005	0.004	0.027	0.009	0.027
Ours-CP-384	0.021	0.040	0.010	0.039	0.006	0.007	0.014	0.015	0.042
Ours-VM-192-SH	0.015	0.030	0.008	0.021	0.003	0.003	0.016	0.016	0.025
Ours-VM-48	0.016	0.031	0.008	0.025	0.004	0.004	0.015	0.013	0.026
Ours-VM-192-15k	0.015	0.033	0.008	0.022	0.004	0.004	0.015	0.011	0.026
Ours-VM-192-30k	0.012	0.024	0.006	0.016	0.003	0.003	0.011	0.009	0.021

Table 9: Quantitative results on each scene from the **Synthetic-NSVF** [18] dataset.

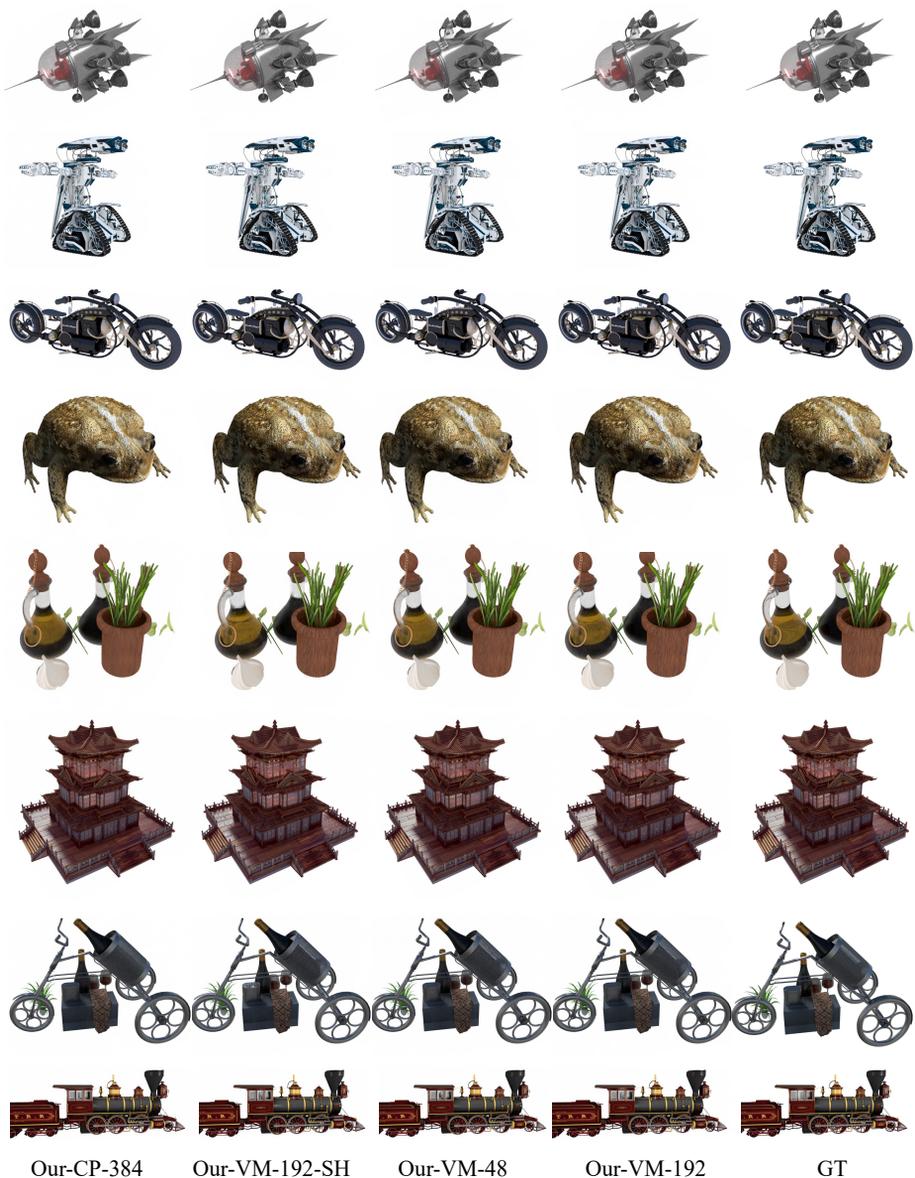


Fig. 7: Our rendering results on NSVF [18] dataset. From top to bottom: Space-ship, Robot, Toad, Lifestyle, Palace, Wineholder, Steamtrain.

Methods	Avg.	<i>Ignatius</i>	<i>Truck</i>	<i>Barn</i>	<i>Caterpillar</i>	<i>Family</i>
PSNR\uparrow						
SRN [36]	24.10	26.70	22.62	22.44	21.14	27.57
NeRF [24]	25.78	25.43	25.36	24.05	23.75	30.29
NSVF [18]	28.48	27.91	26.92	27.16	26.44	33.58
PlenOctrees [47]	27.99	28.19	26.83	26.80	25.29	32.85
Plenoxels [46]	27.43	27.51	26.59	26.07	24.64	32.33
DVG [37]	28.41	28.16	27.15	27.01	26.00	33.75
Ours-CP-384	27.59	27.86	26.25	26.74	24.73	32.39
Ours-VM-192-SH	27.81	27.78	26.73	26.03	25.37	33.12
Ours-VM-48	28.06	28.22	26.81	26.70	25.43	33.12
Ours-VM-192-15k	28.07	28.27	26.57	26.93	25.35	33.22
Ours-VM-192-30k	28.56	28.34	27.14	27.22	26.19	33.92
SSIM\uparrow						
SRN [36]	0.847	0.920	0.832	0.741	0.834	0.908
NeRF [24]	0.864	0.920	0.860	0.750	0.860	0.932
NSVF [18]	0.901	0.930	0.895	0.823	0.900	0.954
PlenOctrees [47]	0.917	0.948	0.914	0.856	0.907	0.962
Plenoxels [46]	0.906	0.943	0.901	0.829	0.902	0.956
DVGO [37]	0.911	0.944	0.906	0.838	0.906	0.962
Ours-CP-384	0.897	0.934	0.885	0.839	0.879	0.948
Ours-VM-192-SH	0.907	0.942	0.900	0.834	0.897	0.960
Ours-VM-48	0.909	0.943	0.902	0.845	0.899	0.957
Ours-VM-192-15k	0.913	0.944	0.905	0.855	0.902	0.960
Ours-VM-192-30k	0.920	0.948	0.914	0.864	0.912	0.965
LPIP$_{VGG} \downarrow$						
PlenOctrees [47]	0.131	0.080	0.130	0.226	0.148	0.069
Plenoxels [46]	0.162	0.102	0.163	0.303	0.166	0.078
DVGO [37]	0.155	0.083	0.160	0.294	0.167	0.070
Ours-CP-384	0.181	0.106	0.202	0.283	0.227	0.088
Ours-VM-192-SH	0.156	0.089	0.161	0.286	0.175	0.069
Ours-VM-48	0.155	0.085	0.161	0.278	0.177	0.074
Ours-VM-192-15k	0.152	0.084	0.162	0.269	0.173	0.071
Ours-VM-192-30k	0.140	0.078	0.145	0.252	0.159	0.064
LPIPS$_{Alex} \downarrow$						
SRN [36]	0.251	0.128	0.266	0.448	0.278	0.134
NeRF [24]	0.198	0.111	0.192	0.395	0.196	0.098
NSVF [18]	0.155	0.106	0.148	0.307	0.141	0.063
DVGO [37]	0.148	0.090	0.145	0.290	0.152	0.064
Ours-CP-384	0.144	0.089	0.154	0.237	0.176	0.063
Ours-VM-192-SH	0.164	0.098	0.168	0.309	0.175	0.072
Ours-VM-48	0.145	0.089	0.145	0.266	0.161	0.066
Ours-VM-192-15k	0.140	0.087	0.150	0.240	0.157	0.066
Ours-VM-192-30k	0.125	0.081	0.129	0.217	0.139	0.057

Table 10: Quantitative results on each scene from the **Tanks&Temples** [15] dataset.

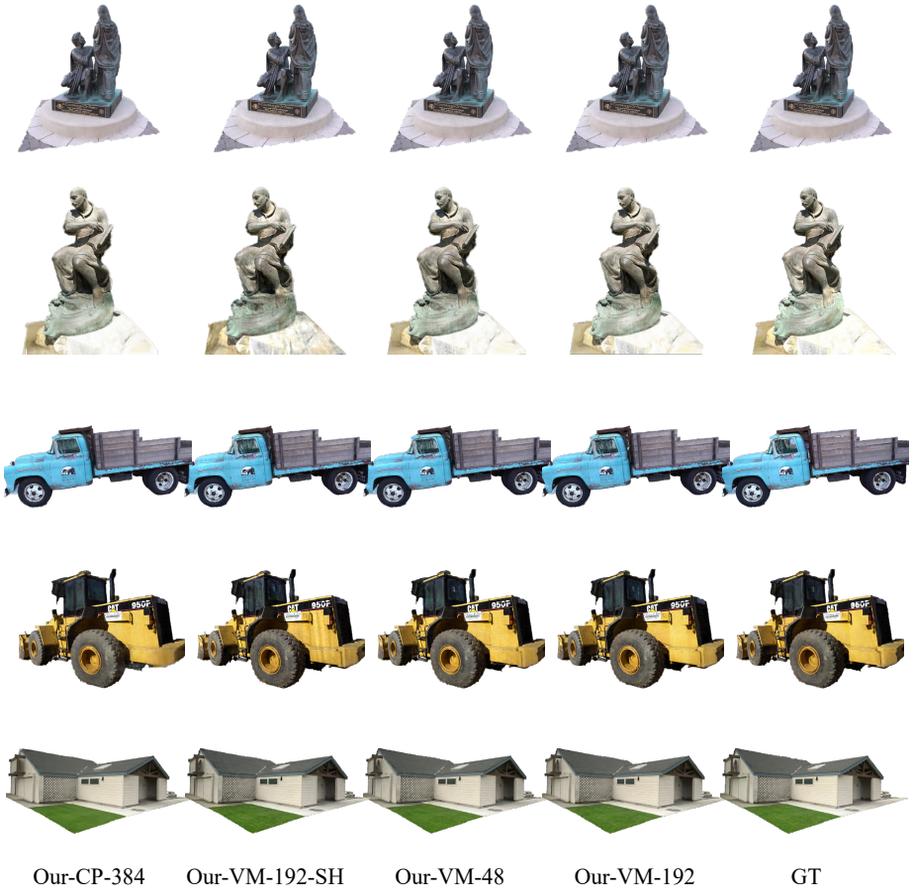
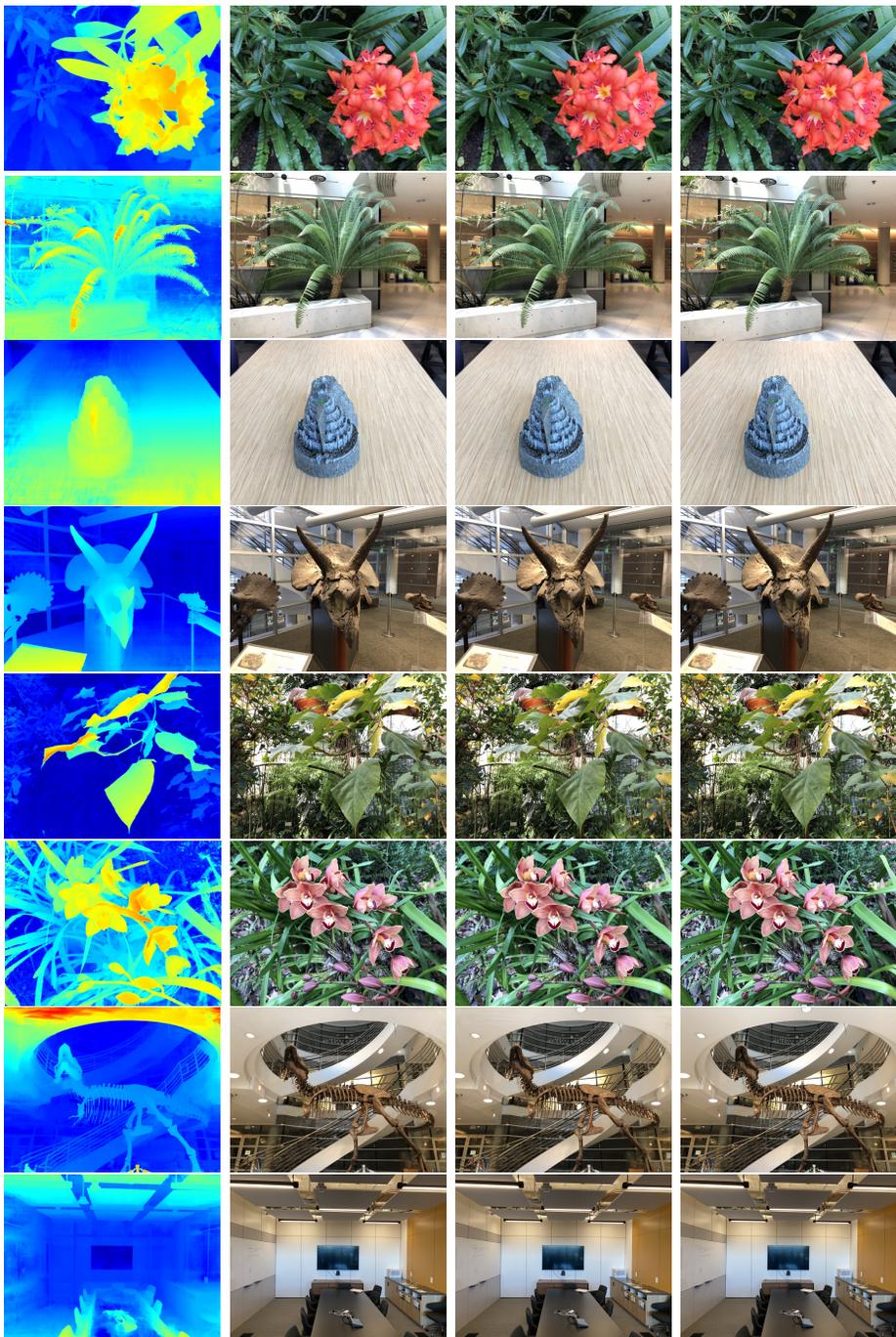


Fig. 8: Our rendering results on **Tanks&Temples**[15] dataset. From top to bottom: Family, Ignatius, Truck, Caterpillar, Barn.

Methods	Avg.	<i>Room</i>	<i>Fern</i>	<i>Leaves</i>	<i>Fortress</i>	<i>Orchids</i>	<i>Flower</i>	<i>T-Rex</i>	<i>Horns</i>
PSNR ↑									
NeRF [24]	26.50	32.70	25.17	20.92	31.16	20.36	27.40	26.80	27.45
Plenoxels [46]	26.29	30.22	25.46	21.41	31.09	20.24	27.83	26.48	27.58
Ours-VM-48	26.51	31.80	25.31	21.34	31.14	20.02	28.22	26.61	27.64
Ours-VM-96	26.73	32.35	25.27	21.30	31.36	19.87	28.60	26.97	28.14
SSIM ↑									
NeRF [24]	0.811	0.948	0.792	0.690	0.881	0.641	0.827	0.880	0.828
Plenoxels [46]	0.839	0.937	0.832	0.760	0.885	0.687	0.862	0.890	0.857
Ours-VM-48	0.832	0.946	0.816	0.746	0.889	0.655	0.859	0.890	0.859
Ours-VM-96	0.839	0.952	0.814	0.752	0.897	0.649	0.871	0.900	0.877
LPIPS _{VGG} ↓									
NeRF [24]	0.250	0.178	0.280	0.316	0.171	0.321	0.219	0.249	0.268
Plenoxels [46]	0.210	0.192	0.224	0.198	0.180	0.242	0.179	0.238	0.231
Ours-VM-48	0.217	0.181	0.237	0.230	0.159	0.283	0.187	0.236	0.221
Ours-VM-96	0.204	0.167	0.237	0.217	0.148	0.278	0.169	0.221	0.196
LPIPS _{Alex} ↓									
Ours-VM-48	0.135	0.093	0.161	0.167	0.084	0.204	0.121	0.108	0.146
Ours-VM-96	0.124	0.082	0.155	0.153	0.075	0.201	0.106	0.099	0.123

Table 11: Quantitative results on each scene from the **forward-facing** [18] dataset.



Our-VM-192

Our-VM-48

Our-VM-192

GT

Fig. 9: Our rendering results on **forward-facing** [18] dataset. From top to bottom: Flower, Fern, Fortress, Horn, Leaves, Orchids, T-Rex, Room.