# FastNeRF: High-Fidelity Neural Rendering at 200FPS

Stephan J. Garbin*        Marek Kowalski*        Matthew Johnson
Jamie Shotton        Julien Valentin

Microsoft

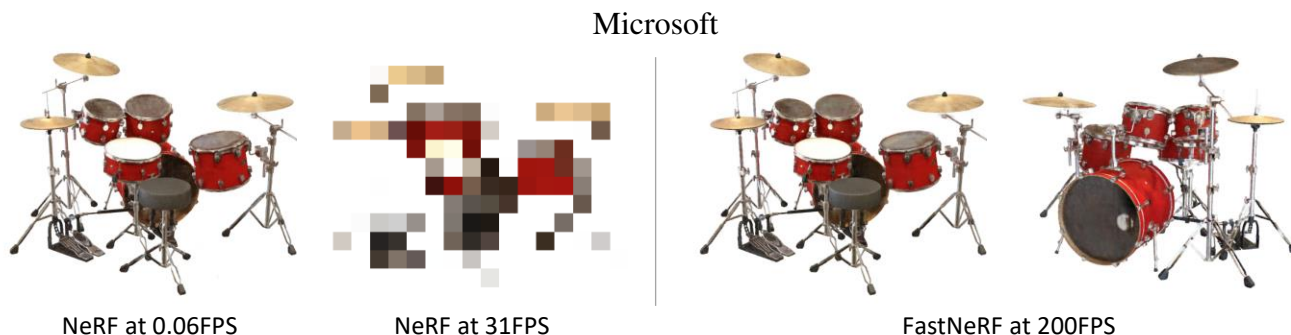NeRF at 0.06FPS          NeRF at 31FPS          FastNeRF at 200FPS

Figure 1. **FastNeRF** renders high-resolution photorealistic novel views of objects at hundreds of frames per second. Comparable existing methods, such as NeRF, are orders of magnitude slower and can only render very low resolution images at interactive rates.

## Abstract

*Recent work on Neural Radiance Fields (NeRF) showed how neural networks can be used to encode complex 3D environments that can be rendered photorealistically from novel viewpoints. Rendering these images is very computationally demanding and recent improvements are still a long way from enabling interactive rates, even on high-end hardware. Motivated by scenarios on mobile and mixed reality devices, we propose FastNeRF, the first NeRF-based system capable of rendering high fidelity photorealistic images at 200Hz on a high-end consumer GPU. The core of our method is a graphics-inspired factorization that allows for (i) compactly caching a deep radiance map at each position in space, (ii) efficiently querying that map using ray directions to estimate the pixel values in the rendered image. Extensive experiments show that the proposed method is 3000 times faster than the original NeRF algorithm and at least an order of magnitude faster than existing work on accelerating NeRF, while maintaining visual quality and extensibility.*

## 1. Introduction

Rendering scenes in real-time at photorealistic quality has long been a goal of computer graphics. Traditional approaches such as rasterization and ray-tracing often require

---

*Denotes equal contribution

Project website: https://microsoft.github.io/FastNeRF

significant manual effort in designing or pre-processing the scene in order to achieve both quality and speed. Recently, neural rendering [11, 29, 19, 27, 21] has offered a disruptive alternative: involve a neural network in the rendering pipeline to output either images directly [26, 20, 29] or to model implicit functions that represent a scene appropriately [5, 25, 31, 27, 42]. Beyond rendering, some of these approaches implicitly reconstruct a scene from static or moving cameras [27, 38, 1], thereby greatly simplifying the traditional reconstruction pipelines used in computer vision.

One of the most prominent recent advances in neural rendering is Neural Radiance Fields (NeRF) [27] which, given a handful of images of a static scene, learns an implicit volumetric representation of the scene that can be rendered from novel viewpoints. The rendered images are of high quality and correctly retain thin structures, view-dependent effects, and partially-transparent surfaces. NeRF has inspired significant follow-up work that has addressed some of its limitations, notably extensions to dynamic scenes [33, 8, 48], relighting [2, 3, 39], and incorporation of uncertainty [24].

One common challenge to all of the NeRF-based approaches is their high computational requirements for rendering images. The core of this challenge resides in NeRF's volumetric scene representation. More than 100 neural network calls are required to render a single image pixel, which translates into several seconds being required to render low-resolution images on high-end GPUs. Recent explorations [37, 17, 30, 18] conducted with the aim of improving NeRF's computational requirements reduced the

render time by up to 50×. While impressive, these advances are still a long way from enabling real-time rendering on consumer-grade hardware. Our work bridges this gap while maintaining quality, thereby opening up a wide range of new applications for neural rendering. Furthermore, our method could form the fundamental building block for neural rendering at high resolutions.

To achieve this goal, we use caching to trade memory for computational efficiency. As NeRF is fundamentally a function of positions $\mathbf{p} \in \mathbb{R}^3$ and ray directions $\mathbf{d} \in \mathbb{R}^2$ to color $\mathbf{c} \in \mathbb{R}^3$ (RGB) and a scalar density $\sigma$, a naïve approach would be to build a cached representation of this function in the space of its input parameters. Since $\sigma$ only depends on $\mathbf{p}$, it can be cached using existing methodologies. The color $\mathbf{c}$, however, is a function of both ray direction $\mathbf{d}$ *and* position $\mathbf{p}$. If this 5 dimensional space were to be discretized with 512 bins per dimension, a cache of around 176 terabytes of memory would be required – dramatically more than is available on current consumer hardware.

Ideally, we would treat directions and positions separately and thus avoid the polynomial explosion in required cache space. Fortunately this problem is not unique to NeRF; the rendering equation (modulo the wavelength of light and time) is also a function of $\mathbb{R}^5$ and solving it efficiently is the primary challenge of real-time computer graphics. As such there is a large body of research which investigates ways of approximating this function as well as efficient means of evaluating the integral. One of the fastest approximations of the rendering equation involves the use of spherical harmonics such that the integral results in a dot product of the harmonic coefficients of the material model and lighting model. Inspired by this efficient approximation, we factorize the problem into two independent functions whose outputs are combined using their inner product to produce RGB values. The first function takes the form of a Multi-Layer Perceptron (MLP) that is conditioned on position in space and returns a compact deep radiance map parameterized by $D$ components. The second function is an MLP conditioned on ray direction that produces $D$ weights for the $D$ deep radiance map components.

This factorized architecture, which we call FastNeRF, allows for independently caching the position-dependent and ray direction-dependent outputs. Assuming that $k$ and $l$ denote the number of bins for positions and ray directions respectively, caching NeRF would have a memory complexity of $\mathcal{O}(k^3 l^2)$. In contrast, caching FastNeRF would have a complexity of $\mathcal{O}(k^3 * (1 + D * 3) + l^2 * D)$. As a result of this reduced memory complexity, FastNeRF can be cached in the memory of a high-end consumer GPU, thus enabling very fast function lookup times that in turn lead to a dramatic increase in test-time performance.

While caching does consume a significant amount of memory, it is worth noting that current implementations of NeRF also have large memory requirements. A single forward pass of NeRF requires performing hundreds of forward passes through an eight layer 256 hidden unit MLP *per pixel*. If pixels are processed in parallel for efficiency this consumes large amounts of memory, even at moderate resolutions. Since many natural scenes (*e.g.* a living room, a garden) are sparse, we are able to store our cache sparsely. In some cases this can make our method actually more memory efficient than NeRF.

In summary, our main contributions are:

- The first NeRF-based system capable of rendering photorealistic novel views at 200FPS, thousands of times faster than NeRF.

- A graphics-inspired factorization that can be compactly cached and subsequently queried to compute the pixel values in the rendered image.

- A blueprint detailing how the proposed factorization can efficiently run on the GPU.

## 2. Related work

FastNeRF belongs to the family of Neural Radiance Fields methods [27] and is trained to learn an implicit, compressed deep radiance map parameterized by position and view direction that provides color and density estimates. Our method differs from [27] in the structure of the implicit model, changing it in such a way that positional and directional components can be discretized and stored in sparse 3D grids.

This also differentiates our method from models that use a discretized grid at training time, such as Neural Volumes [20] or Deep Reflectance Volumes [3]. Due to the memory requirements associated with generating and holding 3D volumes at training time, the output image resolution of [20, 3] is limited by the maximum volume size of $128^3$. In contrast, our method uses volumes as large as $1024^3$.

Subsequent to [20, 3], the problem of parameter estimation for an MLP with low dimensional input coordinates was addressed using Fourier feature encodings. After use in [43], this was popularized in [27, 51] and explored in greater detail by [42].

Neural Radiance Fields (NeRF) [27] first showed convincing compression of a light field utilizing Fourier features. Using an entirely implicit model, NeRF is not bound to any voxelized grid, but only to a specific domain. Despite impressive results, one disadvantage of this method is that a complex MLP has to be called for every sample along each ray, meaning hundreds of MLP invocations per image pixel.

One way to speed up NeRF's inference is to scale to multiple processors. This works particularly well because every pixel in NeRF can be computed independently. For an $800 \times 800$ pixel image, JaxNeRF [6] achieves an inference

speed of 20.77 seconds on an Nvidia Tesla V100 GPU, 2.65 seconds on *8* V100 GPUS, and 0.35 seconds on 128 second generation Tensor Processing Units. Similarly, the volumetric integration domain can be split and separate models used for each part. This approach is taken in the Decomposed Radiance Fields method [37], where an integration domain is divided into Voronoi cells. This yields better test scores, and a speedup of up to $3\times$.

A different way to increase efficiency is to realize that natural scenes tend to be volumetrically sparse. Thus, efficiency can be gained by skipping empty regions of space. This amounts to importance sampling of the occupancy distribution of the integration domain. One approach is to use a voxel grid in combination with the implicit function learned by the MLP, as proposed in Neural Sparse Voxel Fields (NSVFs) [18] and Neural Geometric Level of Detail [40] (for SDFs only), where a dynamically constructed sparse octree is used to represent scene occupancy. As a network still has to be queried inside occupied voxels, however, NSVFs takes between 1 and 4 seconds to render an $800^2$ image, with decreases to PSNR at the lower end of those timings. Our method is orders of magnitude faster in a similar scenario.

Another way to represent the importance distribution is via depth prediction for each pixel. This approach is taken in [30], which is concurrent to ours and achieves roughly 15FPS for $800^2$ images at reduced quality or roughly half that for quality comparable to NeRF.

Orthogonal to this, AutoInt [17] showed that a neural network can be used to approximate the integrals along each ray with far fewer samples. While significantly faster than NeRF, this still does not provide interactive frame rates.

What differentiates our method from those described above is that FastNeRF's proposed decomposition, and subsequent caching, lets us avoid calls to an MLP at inference time entirely. This makes our method faster in absolute terms even on a single machine.

It is worth noting that our method does not address training speed, as for example [41]. In that case, the authors propose improving the training speed of NeRF models by finding initialisation through meta learning.

Finally, there are orthogonal neural rendering strategies capable of fast inference, such as Neural Point-Based Graphics [13], Mixture of Volumetric Primitives [21] or Pulsar [14], which use forms of geometric primitives and rasterization. In this work, we only deal with NeRF-like implicit functional representations paired with ray tracing.

## 3. Method

In this section we describe FastNeRF, a method that is 3000 times faster than the original Neural Radiance Fields (NeRF) system [27] (Section 3.1). This breakthrough allows for rendering high-resolution photorealistic images at

over 200Hz on high-end consumer hardware. The core insight of our approach (Section 3.2) consists of factorizing NeRF into two neural networks: a position-dependent network that produces a deep radiance map and a direction-dependent network that produces weights. The inner product of the weights and the deep radiance map estimates the color in the scene at the specified position and as seen from the specified direction. This architecture, which we call FastNeRF, can be efficiently cached (Section 3.3), significantly improving test time efficiency whilst preserving the visual quality of NeRF. See Figure 2 for a comparison of the NeRF and FastNeRF network architectures.

### 3.1. Neural Radiance Fields

A Neural Radiance Field (NeRF) [27] captures a volumetric 3D representation of a scene within the weights of a neural network. NeRF's neural network $\mathcal{F}_{NeRF} : (\boldsymbol{p}, \boldsymbol{d}) \mapsto (\boldsymbol{c}, \sigma)$ maps a 3D position $\boldsymbol{p} \in \mathbb{R}^3$ and a ray direction $\boldsymbol{d} \in \mathbb{R}^2$ to a color value $\boldsymbol{c}$ and transparency $\sigma$. See the left side of Figure 2 for a diagram of $\mathcal{F}_{NeRF}$.

In order to render a single image pixel, a ray is cast from the camera center, passing through that pixel and into the scene. We denote the direction of this ray as $\boldsymbol{d}$. A number of 3D positions $(\boldsymbol{p}_1, \cdots, \boldsymbol{p}_N)$ are then sampled along the ray between its near and far bounds defined by the camera parameters. The neural network $\mathcal{F}_{NeRF}$ is evaluated at each position $\boldsymbol{p}_i$ and ray direction $\boldsymbol{d}$ to produce color $\boldsymbol{c}_i$ and transparency $\sigma_i$. These intermediate outputs are then integrated as follows to produce the final pixel color $\hat{\boldsymbol{c}}$:

$$\hat{\boldsymbol{c}} = \sum_{i=1}^{N} T_i(1 - \exp(-\sigma_i \delta_i))\boldsymbol{c}_i, \qquad (1)$$

where $T_i = \exp(-\sum_{j=i}^{i-1} \sigma_j \delta_j)$ is the transmittance and $\delta_i = (\boldsymbol{p}_{i+1} - \boldsymbol{p}_i)$ is the distance between the samples. Since $\mathcal{F}_{NeRF}$ depends on ray directions, NeRF has the ability to model viewpoint-dependent effects such as specular reflections, which is one key dimension in which NeRF improves upon traditional 3D reconstruction methods.

Training a NeRF network requires a set of images of a scene as well as the extrinsic and intrinsic parameters of the cameras that captured the images. In each training iteration, a subset of pixels from the training images are chosen at random and for each pixel a 3D ray is generated. Then, a set of samples is selected along each ray and the pixel color $\hat{\boldsymbol{c}}$ is computed using Equation (1). The training loss is the mean squared difference between $\hat{\boldsymbol{c}}$ and the ground truth pixel value. For further details please refer to [27].

While NeRF renders photorealistic images, it requires calling $\mathcal{F}_{NeRF}$ a large number of times to produce a single image. With the default number of samples per pixel $N = 192$ proposed in NeRF [27], nearly 400 million calls to $\mathcal{F}_{NeRF}$ are required to compute a single high definition (1080p) image. Moreover the intermediate outputs of
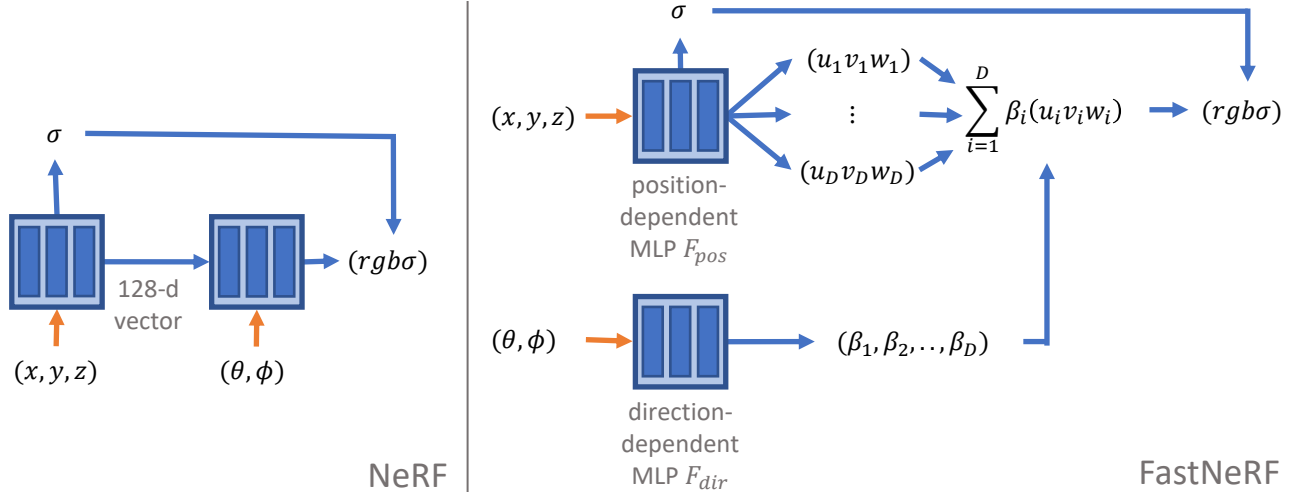
Figure 2. Left: NeRF neural network architecture. $(x, y, z)$ denotes the input sample position, $(\theta, \phi)$ denotes the ray direction and $(r, g, b, \sigma)$ are the output color and transparency values. Right: our FastNeRF architecture splits the same task into two neural networks that are amenable to caching. The position-dependent network $\mathcal{F}_{pos}$ outputs a deep radiance map $(\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w})$ consisting of $D$ components, while the $\mathcal{F}_{dir}$ outputs the weights for those components $(\beta_1, \ldots, \beta_D)$ given a ray direction as input.

this process would take hundreds of gigabytes of memory. Even on high-end consumer GPUs, this constrains the original method to be executed over several batches even for medium resolution ($800 \times 800$) images, leading to additional computational overheads.

### 3.2. Factorized Neural Radiance Fields

Taking a step away from neural rendering for a moment, we recall that in traditional computer graphics, the rendering equation [10] is an integral of the form

$$L_o(\boldsymbol{p}, \boldsymbol{d}) = \int_{\Omega} f_r(\boldsymbol{p}, \boldsymbol{d}, \boldsymbol{\omega_i}) L_i(\boldsymbol{p}, \boldsymbol{\omega_i})(\boldsymbol{\omega_i} \cdot \boldsymbol{n}) d\boldsymbol{\omega_i}, \quad (2)$$

where $L_o(\boldsymbol{p}, \boldsymbol{d})$ is the radiance leaving the point $\boldsymbol{p}$ in direction $\boldsymbol{d}$, $f_r(\boldsymbol{p}, \boldsymbol{d}, \boldsymbol{\omega_i})$ is the reflectance function capturing the material properties at position $\boldsymbol{p}$, $L_i(\boldsymbol{p}, \boldsymbol{\omega_i})$ describes the amount of light reaching $\boldsymbol{p}$ from direction $\boldsymbol{\omega_i}$, and $\boldsymbol{n}$ corresponds to the direction of the surface normal at $\boldsymbol{p}$. Given its practical importance, evaluating this integral efficiently has been a subject of active research for over three decades [10, 44, 9, 35]. One efficient way to evaluate the rendering equation is to approximate $f_r(\boldsymbol{p}, \boldsymbol{d}, \boldsymbol{\omega_i})$ and $L_i(\boldsymbol{p}, \boldsymbol{\omega_i})$ using spherical harmonics [4, 47]. In this case, evaluating the integral boils down to a dot product between the coefficients of both spherical harmonics approximations.

With this insight in mind, we can return to neural rendering. In the case of NeRF, $\mathcal{F}_{NeRF} : (\boldsymbol{p}, \boldsymbol{d})$ can also be interpreted as returning the radiance leaving point $\boldsymbol{p}$ in direction $\boldsymbol{d}$. This key insight leads us to propose a new network architecture for NeRF, which we call FastNeRF. This novel architecture consists in splitting NeRF's neural network $\mathcal{F}_{NeRF}$ into two networks: one that depends only

on positions $\boldsymbol{p}$ and one that depends only on the ray directions $\boldsymbol{d}$. Similarly to evaluating the rendering equation using spherical harmonics, our position-dependent $\mathcal{F}_{pos}$ and direction-dependent $\mathcal{F}_{dir}$ functions produce outputs that are combined using the dot product to obtain the color value at position $\boldsymbol{p}$ observed from direction $\boldsymbol{d}$. Crucially, this factorization splits a single function that takes inputs in $\mathbb{R}^5$ into two functions that take inputs in $\mathbb{R}^3$ and $\mathbb{R}^2$. As explained in the following section, this makes caching the outputs of the network possible and allows accelerating NeRF by 3 orders of magnitude on consumer hardware. See Figure 3 for a visual representation of the achieved speedup.

The position-dependent and direction-dependent functions of FastNeRF are defined as follows:

$$\mathcal{F}_{\text{pos}} : \boldsymbol{p} \mapsto \{\sigma, (\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w})\}, \quad (3)$$

$$\mathcal{F}_{\text{dir}} : \boldsymbol{d} \mapsto \boldsymbol{\beta}, \quad (4)$$

where $\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w}$ are $D$-dimensional vectors that form a deep radiance map describing the view-dependent radiance at position $\boldsymbol{p}$. The output of $\mathcal{F}_{dir}$, $\boldsymbol{\beta}$, is a $D$-dimensional vector of weights for the $D$ components of the deep radiance map. The inner product of the weights and the deep radiance map

$$\boldsymbol{c} = (r, g, b) = \sum_{i=1}^{D} \beta_i(u_i, v_i, w_i) = \boldsymbol{\beta}^T \cdot (\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w}) \quad (5)$$

results in the estimated color $\boldsymbol{c} = (r, g, b)$ at position $\boldsymbol{p}$ observed from direction $\boldsymbol{d}$.

When the two functions are combined using Equation (5), the resulting function $\mathcal{F}_{\text{FastNeRF}}(\boldsymbol{p}, \boldsymbol{d}) \mapsto (\boldsymbol{c}, \sigma)$ has the same signature as the function $\mathcal{F}_{NeRF}$ used in NeRF. This effectively allows for using the FastNeRF architecture as a drop-in replacement for NeRF's architecture at runtime.

4

## 3.3. Caching

Given the large number of samples that need to be evaluated to render a single pixel, the cost of computing $\mathcal{F}$ dominates the total cost of NeRF rendering. Thus, to accelerate NeRF, one could attempt to reduce the test-time cost of $\mathcal{F}$ by caching its outputs for a set of inputs that cover the space of the scene. The cache can then be evaluated at a fraction of the time it takes to compute $\mathcal{F}$.

For a trained NeRF model, we can define a bounding box $\mathcal{V}$ that covers the entire scene captured by NeRF. We can then uniformly sample $k$ values for each of the 3 world-space coordinates $(x, y, z) = \boldsymbol{p}$ within the bounds of $\mathcal{V}$. Similarly, we uniformly sample $l$ values for each of the ray direction coordinates $(\theta, \phi) = \boldsymbol{d}$ with $\theta \in \langle 0, \pi \rangle$ and $\phi \in \langle 0, 2\pi \rangle$. The cache is then generated by computing $\mathcal{F}$ for each combination of sampled $\boldsymbol{p}$ and $\boldsymbol{d}$.

The size of such a cache for a *standard* NeRF model with $k = l = 1024$ and densely stored 16-bit floating point values is approximately 5600 Terabytes. Even for highly sparse volumes, where one would only need to keep 1% of the outputs, the size of the cache would still severely exceed the memory capacity of consumer-grade hardware. This huge memory requirement is caused by the usage of both $\boldsymbol{d}$ and $\boldsymbol{p}$ as input to $\mathcal{F}_{NeRF}$. A separate output needs to be saved for each combination of $\boldsymbol{d}$ and $\boldsymbol{p}$, resulting in a memory complexity in the order of $\mathcal{O}(k^3 l^2)$.

Our FastNeRF architecture makes caching feasible. For $k = l = 1024, D = 8$ the size of two dense caches holding $\{\sigma, (\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w}))\}$ and $\boldsymbol{\beta}$ would be approximately 54 GB. For moderately sparse volumes, where 30% of space is occupied, the memory requirement is low enough to fit into either the CPU or GPU memory of consumer-grade machines. In practice, the choice $k$ and $l$ depends on the scene size and the expected image resolution. For many scenarios a smaller cache of $k = 512, l = 256$ is sufficient, lowering the memory requirements further. Please see the supplementary materials for formulas used to calculate cache sizes for both network architectures.

## 4. Implementation

Aside from using the proposed FastNeRF architecture described in Section 3.2, training FastNeRF is identical to training NeRF [27]. We model $\mathcal{F}_{pos}$ and $\mathcal{F}_{view}$ of FastNeRF using 8 and 4 layer MLPs respectively, with positional encoding applied to inputs [27]. The 8-layer MLP has 384 hidden units, and we use 128 for the view-predicting 4-layer network. The only exception to this is the *ficus* scene, where we found a 256 MLP gave better results for the 8-layer MLP. While this makes the network larger than the baseline, cache lookup speed is not influenced by MLP complexity. Similarly to NeRF, we parameterize the view direction as a 3-dimensional unit vector. Other training-time features
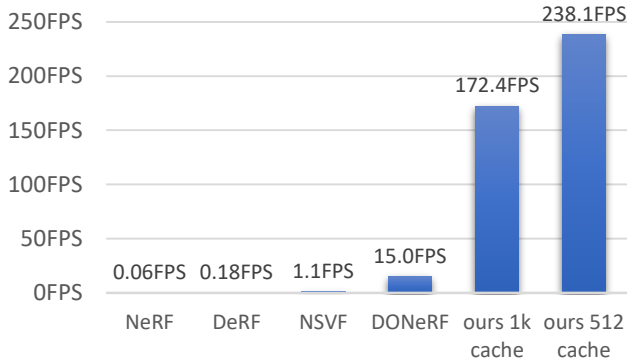


Figure 3. Speed evaluation of our method and prior work [27, 37, 18, 30] on the Lego scene from the Realistic 360 Synthetic [27] dataset, rendered at $800 \times 800$ pixels. For previous methods, when numbers for the Lego scene were not available, we used an optimistic approximation.

and settings including the coarse/fine networks and sample count $N$ follow the original NeRF work [27] (please see supplementary materials).

At test time, both our method and NeRF take a set of camera parameters as input which are used to generate a ray for each pixel in the output. A number of samples are then produced along each ray and integrated following Section 3.1. While FastNeRF can be executed using its neural network representation, its performance is greatly improved when cached. To further improve performance, we use hardware accelerated ray tracing [34] to skip empty space, starting to integrate points along the ray only after the first hit with a collision mesh derived from the density volume, and visiting every voxel thereafter until a ray's transmittance is saturated. The collision mesh can be computed from a sign-distance function derived from the density volume using Marching Cubes [22]. For volumes greater than $512^3$, we first downsample the volume by a factor of two to reduce mesh complexity. We use the same meshing and integration parameters for all scenes and datasets. Value lookup uses nearest neighbour interpolation for $\mathcal{F}_{pos}$ and trilinear sampling for $\mathcal{F}_{view}$, the latter of which can be processed with a Gaussian or other kernel to 'smooth' or edit directional effects. We note that the meshes and volumes derived from the cache can also be used to provide approximations to shadows and global illumination in a path-tracing framework.

## 5. Experiments

We evaluate our method quantitatively and qualitatively on the Realistic 360 Synthetic [27] and Local Light Field Fusion (LLFF) [26] (with additions from [27]) datasets used in the original NeRF paper. While the NeRF synthetic dataset consists of 360 degree views of complex objects, the LLFF dataset consist of forward-facing scenes, with fewer
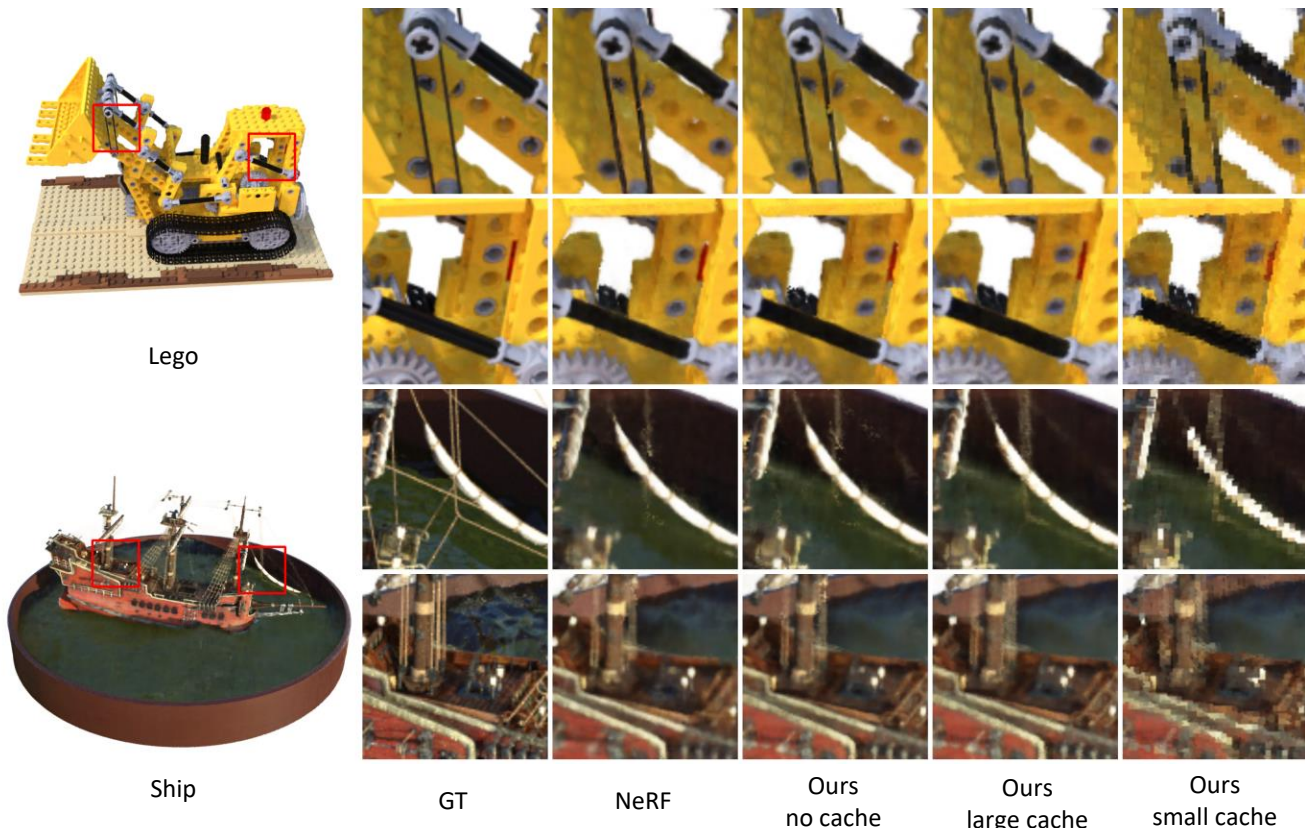
Figure 4. Qualitative comparison of our method vs NeRF on the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$. Varying the cache size allows for trading compute and memory for image quality resembling levels of detail (LOD) in traditional computer graphics.

Table 1. We compare NeRF to our method when not cached to a grid, and when cached at a high resolution in terms of PSNR, SSIM and LPIPS, and also provide the average speed in ms of our method when cached. LLFF denotes the dataset of [26] at $504 \times 378$ pixels, the other dataset is that of [27] at $800^2$ pixels. We use 8 components and a $1024^3$ cache for the NeRF Synthetic 360 scenes, and 6 components at $768^3$ for LLFF. Please see the supplementary material for more detailed results.

| Scene | NeRF | | | Ours - No Cache | | | Ours - Cache | | | Speed |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | |
| **Nerf Synthetic** | 29.54dB | 0.94 | 0.05 | 29.155dB | 0.936 | 0.053 | 29.97dB | 0.941 | 0.053 | 4.2ms |
| **LLFF** | 27.72dB | 0.88 | 0.07 | 27.958dB | 0.888 | 0.063 | 26.035dB | 0.856 | 0.085 | 1.4ms |

images. In all comparisons with NeRF we use the same training parameters as described in the original paper [27].

To assess the rendering quality of FastNeRF we compare its outputs to GT using Peak Signal to Noise Ratio (PSNR), Structural Similarity (SSIM) [45] and perceptual LPIPS [50]. We use the same metrics in our ablation study, where we evaluate the effects of changing the number of components $D$ and the size of the cache. All speed comparisons are run on one machine that is equipped with an Nvidia RTX 3090 GPU.

We keep the size of the view-dependent cache $\mathcal{F}_{\text{view}}^{\text{cache}}$ the same for all experiments at a base resolution of $384^3$. It is the base resolution of the RGBD cache $\mathcal{F}_{\text{pos}}^{\text{cache}}$ that

represents the main trade-off in complexity of runtime and memory vs image quality. Across all results, grid resolution refers to the number of divisions of the longest side of the bounding volume surrounding the object of interest.

**Rendering quality:** Because we adopt the inputs, outputs, and training of NeRF, we retain the compatibility with ray-tracing [46, 10] and the ability to only specify loose volume bounds. At the same time, our factorization and caching do not affect the character of the rendered images to a large degree. Please see Figure 4 and Figure 5 for qualitative comparisons, and Table 1 for a quantitative evaluation. Note that some aliasing artefacts can appear since our model is cached to a grid *after* training. This explains the decrease
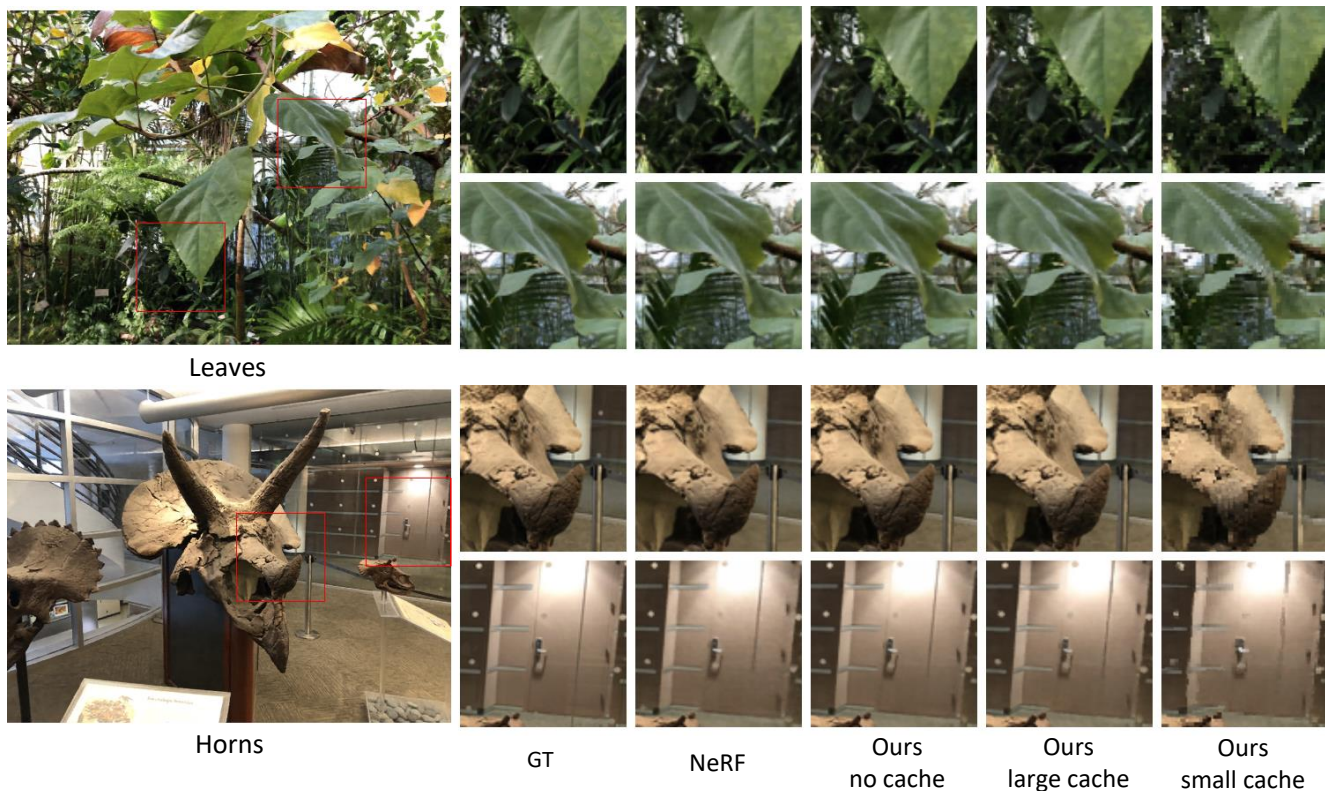
Figure 5. Qualitative comparison of our method vs NeRF on the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.

Table 2. Speed comparison of our method vs NeRF in ms. The *Chair* and *Lego* scenes are at rendered at $800^2$ resolution, *Horns* and *Leaves* scenes at $504 \times 378$. Our method never drops below 100FPS when cached, and is often significantly faster. Please note that our method is slower when not cached due to using larger MLPs; it performs identically to NeRF when using 256 hidden units. We do not compute the highest resolution cache for LLFF scenes because these are less sparse due to the use of normalized device coordinates.

| Scene | NeRF | Ours - No Cache | $256^3$ | $384^3$ | $512^3$ | $768^3$ | $1024^3$ | Speedup over NeRF |
|---|---|---|---|---|---|---|---|---|
| **Chair** | 17.5**K** | 28.2**K** | 0.8 | 1.1 | 1.4 | 2.0 | 2.7 | $6468\times$ - $21828\times$ |
| **Lego** | 17.5**K** | 28.2**K** | 1.5 | 2.1 | 2.8 | 4.2 | 5.6 | $3118\times$ - $11639\times$ |
| **Horns\*** | 3.8**K** | 6.2**K** | 0.5 | 0.7 | 0.9 | 1.2 | - | $3183\times$ - $7640\times$ |
| **Leaves\*** | 3.9**K** | 6.3**K** | 0.6 | 0.8 | 1.0 | 1.5 | - | $2626\times$ - $6566\times$ |

across all metrics as a function of the grid resolution as seen in Table 1. However, both NeRF and FastNeRF mitigate multi-view inconsistencies, ghosting and other such artifacts from multi-view reconstruction. Table 1 demonstrates that at high enough cache resolutions, our method is capable of the same visual quality as NeRF. Figure 4 and Figure 5 further show that smaller caches appear 'pixelated' but retain the overall visual characteristics of the scene. This is similar to how different levels of detail are used in computer graphics [23], and an important innovation on the path towards neurally rendered worlds.

**Cache Resolution:** As shown in Table 1, our method

matches or outperforms NeRF on the dataset of synthetic objects at $1024^3$ cache resolution. Because our method is fast enough to visit every voxel (as opposed to the fixed sample count used in NeRF), it can sometimes achieve better scores by not missing any detail. We observe that a cache of $512^3$ is a good trade-off between perceptual quality, memory and rendering speed for the synthetic dataset. For the LLFF dataset, we found a $768^3$ cache to work best. For the *Ship* scene, an ablation over the grid size, memory, and PSNR is shown in Table 3.

For the view-filling LLFF scenes at $504 \times 378$ pixels, our method sees a slight decrease in metrics when cached

Table 3. Influence of the number of components and grid resolution on PSNR and memory required for caching the ship scene. Note how more factors increase grid sparsity. We find that 8 or 6 components are a reasonable compromise in practice.

| Factors | No Cache | | $256^3$ | | $384^3$ | | $512^3$ | | $768^3$ | |
|---------|----------|--------|---------|--------|---------|--------|---------|--------|---------|--------|
| | PSNR↑ | Memory | PSNR↑ | Memory | PSNR↑ | Memory | PSNR↑ | Memory | PSNR↑ | Memory |
| 4 | 27.11dB | - | 24.81dB | 0.34Gb | 26.29dB | 0.61Gb | 26.94dB | 1.09Gb | 27.54dB | 2.51Gb |
| 6 | 27.12dB | - | 24.82dB | 0.5Gb | 26.34dB | 0.93Gb | 27.0dB | 1.67Gb | 27.58dB | 4.1Gb |
| 8 | 27.24dB | - | 24.89dB | 0.71Gb | 26.42dB | 1.41Gb | 27.1dB | 2.7Gb | 27.72dB | 7.15Gb |
| 16 | 27.68dB | - | 25.07dB | 1.2Gb | 26.77dB | 2.08Gb | 27.55dB | 3.72Gb | 28.3dB | 9.16Gb |



Figure 6. Face images rendered using FastNeRF combined with a deformation field network [32]. Thanks to the use of FastNeRF, expression-conditioned images can be rendered at 30FPS.

at $768^3$, but still produces qualitatively compelling results as shown in Figure 5, where intricate detail is clearly preserved.

**Rendering Speed:** When using a grid resolution of $768^3$, FastNeRF is on average more than $3000\times$ faster than NeRF at largely the same perceptual quality. See Table 2 for a breakdown of run-times across several resolutions of the cache and Figure 3 for a comparison to other NeRF extensions in terms of speed. Note that we log the time it takes our method to fill an RGBA buffer of the image size with the final pixel values, whereas the baseline implementation needs to perform various steps to reshape samples back into images. Our CUDA kernels are not highly optimized, favoring flexibility over maximum performance. It is reasonable to assume that advanced code optimization and further compression of the grid values could lead to further reductions in compute time.

**Number of Components:** While we can see a theoretical improvement of roughly 0.5dB going from 8 to 16 components, we find that 6 or 8 components are sufficient for most scenes. As shown in Table 3, our cache can bring out fine details that are lost when only a fixed sample count is used, compensating for the difference. Table 3 also shows that more components tend to be increasingly sparse when cached, which compensates somewhat for the additional memory usage. Please see the supplementary material for more detailed results on other scenes.

## 6. Application

We demonstrate a proof-of-concept application of FastNeRF for a telepresence scenario by first gathering a dataset of a person performing facial expressions for about 20s in a multi-camera rig consisting of 32 calibrated cameras. We then fit a 3D face model similar to FLAME [15] to obtain expression parameters for each frame. Since the captured scene is not static, we train a FastNeRF model of the scene jointly with a deformation model [32] conditioned on the expression data. The deformation model takes the samples $p$ as input and outputs their updated positions that live in a canonical frame of reference modeled by FastNeRF.

We show example outputs produced using this approach in Figure 6. This method allows us to render $300\times300$ pixel images of the face at 30 fps on a single Nvidia Tesla V100 GPU – around 50 times faster than a setup with a NeRF model. At the resolution we achieve, the face expression is clearly visible and the high frame rate allows for real-time rendering, opening the door to telepresence scenarios. The main limitation in terms of speed and image resolution is the deformation model, which needs to be executed for a large number of samples. We employ a simple pruning method, detailed in the supplementary material, to reduce the number of samples.

Finally, note that while we train this proof-of-concept method on a dataset captured using multiple cameras, the approach can be extended to use only a single camera, similarly to [8].

Since FastNeRF uses the same set of inputs and outputs as those used in NeRF, it is applicable to many of NeRF's extensions. This allows for accelerating existing approaches for: reconstruction of dynamic scenes [32, 36, 16], single-image reconstruction [41], quality improvement [37, 49] and others [30]. With minor modifications, FastNeRF can be applied to even more methods allowing for control over illumination [2] and incorporation of uncertainty [24].

## 7. Conclusion

In this paper we presented FastNeRF, a novel extension to NeRF that enables the rendering of photorealistic images at 200Hz and more on consumer hardware. We achieve significant speedups over NeRF and competing methods by factorizing its function approximator, enabling a caching step that makes rendering memory-bound instead of compute-bound. This allows for the application of NeRF to real-time scenarios.

# References

[1] Mojtaba Bemana, Karol Myszkowski, Hans-Peter Seidel, and Tobias Ritschel. X-fields: Implicit neural view-, light- and time-image interpolation. *ACM Transactions on Graphics (Proc. SIGGRAPH Asia 2020)*, 39(6), 2020. 1

[2] Sai Bi, Zexiang Xu, Pratul Srinivasan, Ben Mildenhall, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Neural reflectance fields for appearance acquisition, 2020. 1, 8

[3] Sai Bi, Zexiang Xu, Kalyan Sunkavalli, Miloš Hašan, Yannick Hold-Geoffroy, David Kriegman, and Ravi Ramamoorthi. Deep reflectance volumes: Relightable reconstructions from multi-view photometric images, 2020. 1, 2

[4] Brian Cabral, Nelson Max, and Rebecca Springmeyer. Bidirectional reflection functions from surface bump maps. 21(4):273–281, Aug. 1987. 4

[5] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. *CoRR*, abs/1812.02822, 2018. 1

[6] Boyang Deng, Jonathan T. Barron, and Pratul P. Srinivasan. JaxNeRF: an efficient JAX implementation of NeRF, 2020. 2

[7] Julian Fong, Magnus Wrenninge, Christopher Kulla, and Ralf Habel. Production volume rendering: Siggraph 2017 course. In *ACM SIGGRAPH 2017 Courses*, SIGGRAPH '17, New York, NY, USA, 2017. Association for Computing Machinery. 12

[8] Guy Gafni, Justus Thies, Michael Zollhöfer, and Matthias Nießner. Dynamic neural radiance fields for monocular 4d facial avatar reconstruction. *arXiv preprint arXiv:2012.03065*, 2020. 1, 8

[9] Iliyan Georgiev, Jaroslav Křivánek, Tomáš Davidovič, and Philipp Slusallek. Light transport simulation with vertex connection and merging. *ACM Trans. Graph.*, 31(6):192:1–192:10, Nov. 2012. 4

[10] James T. Kajiya. The rendering equation. *SIGGRAPH Comput. Graph.*, 20(4):143–150, Aug. 1986. 4, 6

[11] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks, 2019. 1

[12] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. 11

[13] Maria Kolos, Artem Sevastopolsky, and Victor Lempitsky. Transpr: Transparency ray-accumulating neural 3d scene point renderer, 2020. 3

[14] Christoph Lassner and Michael Zollhöfer. Pulsar: Efficient sphere-based neural rendering, 2020. 3

[15] Tianye Li, Timo Bolkart, Michael J Black, Hao Li, and Javier Romero. Learning a model of facial shape and expression from 4d scans. *ACM Trans. Graph.*, 36(6):194–1, 2017. 8

[16] Zhengqi Li, Simon Niklaus, Noah Snavely, and Oliver Wang. Neural scene flow fields for space-time view synthesis of dynamic scenes. *arXiv preprint arXiv:2011.13084*, 2020. 8

[17] David B Lindell, Julien NP Martel, and Gordon Wetzstein. Autoint: Automatic integration for fast neural volume rendering. *arXiv preprint arXiv:2012.01714*, 2020. 1, 3

[18] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields, 2021. 1, 3, 5

[19] Stephen Lombardi, Jason M. Saragih, Tomas Simon, and Yaser Sheikh. Deep appearance models for face rendering. *CoRR*, abs/1808.00362, 2018. 1

[20] 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. 1, 2

[21] Stephen Lombardi, Tomas Simon, Gabriel Schwartz, Michael Zollhoefer, Yaser Sheikh, and Jason Saragih. Mixture of volumetric primitives for efficient neural rendering, 2021. 1, 3

[22] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. *SIGGRAPH Comput. Graph.*, 21(4):163–169, Aug. 1987. 5

[23] David Luebke, Martin Reddy, Jonathan D. Cohen, Amitabh Varshney, Benjamin Watson, and Robert Huebner. *Level of Detail for 3D Graphics*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 2002. 7

[24] 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, 2021. 1, 8

[25] Lars M. Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. *CoRR*, abs/1812.03828, 2018. 1

[26] 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), July 2019. 1, 5, 6, 7, 12, 17, 18, 19, 20

[27] 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. *arXiv preprint arXiv:2003.08934*, 2020. 1, 2, 3, 5, 6, 11, 12, 15, 16, 17

[28] Ken Museth. Hierarchical digital differential analyzer for efficient ray-marching in openvdb. In *ACM SIGGRAPH 2014 Talks*, SIGGRAPH '14, New York, NY, USA, 2014. Association for Computing Machinery. 11

[29] Oliver Nalbach, Elena Arabadzhiyska, Dushyant Mehta, H-P Seidel, and Tobias Ritschel. Deep shading: convolutional neural networks for screen space shading. In *Computer graphics forum*, volume 36, pages 65–78. Wiley Online Library, 2017. 1

[30] Thomas Neff, Pascal Stadlbauer, Mathias Parger, Andreas Kurz, Chakravarty R. Alla Chaitanya, Anton Kaplanyan, and Markus Steinberger. Donerf: Towards real-time rendering of neural radiance fields using depth oracle networks, 2021. 1, 3, 5, 8

[31] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. *CoRR*, abs/1901.05103, 2019. 1

[32] Keunhong Park, Utkarsh Sinha, Jonathan T Barron, Sofien Bouaziz, Dan B Goldman, Steven M Seitz, and Ricardo-Martin Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 8, 14

[33] Keunhong Park, Utkarsh Sinha, Jonathan T. Barron, Sofien Bouaziz, Dan B Goldman, Steven M. Seitz, and Ricardo Martin-Brualla. Deformable neural radiance fields. *arXiv preprint arXiv:2011.12948*, 2020. 1

[34] Steven G Parker, James Bigler, Andreas Dietrich, Heiko Friedrich, Jared Hoberock, David Luebke, David McAllister, Morgan McGuire, Keith Morley, Austin Robison, et al. Optix: a general purpose ray tracing engine. *Acm transactions on graphics (tog)*, 29(4):1–13, 2010. 5

[35] M. Pharr, W. Jakob, and G. Humphreys. *Physically Based Rendering: From Theory to Implementation*. Elsevier Science, 2016. 4

[36] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. *arXiv preprint arXiv:2011.13961*, 2020. 8

[37] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields. *arXiv preprint arXiv:2011.12490*, 2020. 1, 3, 5, 8

[38] Gernot Riegler and Vladlen Koltun. Stable view synthesis, 2020. 1

[39] Pratul P. Srinivasan, Boyang Deng, Xiuming Zhang, Matthew Tancik, Ben Mildenhall, and Jonathan T. Barron. Nerv: Neural reflectance and visibility fields for relighting and view synthesis, 2020. 1

[40] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes, 2021. 3

[41] 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, 2020. 3, 8

[42] Matthew Tancik, Pratul P. Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan T. Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains, 2020. 1, 2

[43] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 2

[44] Eric Veach and Leonidas J. Guibas. Optimally combining sampling techniques for monte carlo rendering. In *Proceedings of the 22nd Annual Conference on Computer Graphics and Interactive Techniques*, SIGGRAPH '95, page 419–428, New York, NY, USA, 1995. Association for Computing Machinery. 4

[45] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. 6

[46] Turner Whitted. An improved illumination model for shaded display. In *ACM Siggraph 2005 Courses*, pages 4–es. 2005. 6

[47] Lifan Wu, Guangyan Cai, Shuang Zhao, and Ravi Ramamoorthi. Analytic spherical harmonic gradients for realtime rendering with many polygonal area lights. *ACM Trans. Graph.*, 39(4), July 2020. 4

[48] Wenqi Xian, Jia-Bin Huang, Johannes Kopf, and Changil Kim. Space-time neural irradiance fields for free-viewpoint video. *arXiv preprint arXiv:2011.12950*, 2020. 1

[49] Kai Zhang, Gernot Riegler, Noah Snavely, and Vladlen Koltun. Nerf++: Analyzing and improving neural radiance fields. *arXiv preprint arXiv:2010.07492*, 2020. 8

[50] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. *CoRR*, abs/1801.03924, 2018. 6

[51] Ellen D. Zhong, Tristan Bepler, Joseph H. Davis, and Bonnie Berger. Reconstructing continuous distributions of 3d protein structure from cryo-em images, 2020. 2

## 8. Overview

We use the supplementary material to provide more detailed results and algorithms. We urge the reader to see our video, in which we show results for all 16 scenes we tested on along with providing an intuition for our method. On a practical note, we also provide guidance on getting smooth results for anyone adapting our method to their own work.

Please note that parameters chosen throughout this work for the rendering algorithm favour the highest possible quality. It is possible to significantly increase rendering speeds by sacrificing a small amount of visual quality for real-world applications.

### 8.1. 'Motion' Smoothness

We notice that our method can exhibit more flickering artefacts on the NeRF 360 Synthetic dataset than the baseline *when trained using the same settings as NeRF*, which we find only noticeable in motion. These direction-dependent artefacts are due to overfitting of the view-dependent MLP and can be successfully mitigated in one of two ways (which can also be combined). Sticking with identical parameters as NeRF, we can simply choose a resolution between $16^3$ and $64^3$ for the direction-dependent cache, which uses interpolation at lookup by default. Alternatively, we can adjust the Fourier feature encoding of the view directions at training time. In [27], the Fourier encoding maps from $\mathbb{R}$ to $\mathbb{R}^{2L}$, where $L$ is 10 for position, and 4 for direction. For FastNeRF, setting $L = 1$ actually works best for the NeRF 360 Synthetic dataset. Please note that we do not use these approaches when computing metrics in the text to ensure the same settings are used for all datasets, but that using them would slightly improve the results on synthetic data.

**Effect of modified Fourier feature encoding on metrics:** Using an 8 layer 256 hidden unit MLP for position, and a 4 layer 128 MLP for direction, we can obtain an average PSNR of $29.748dB$ for $L = 1$, $29.646dB$ for $L = 2$, $29.449dB$ for $L = 3$ on synthetic data, with smoothness of results being reflected in these numbers.

**Effect of small direction cache on metrics:** Using a smaller direction cache has a negligible impact on metrics, occasionally actually improving them. This shows that the direction-dependent effects required by FastNeRF are low in frequency. For the scenes where we found we needed to use a smaller cache for moving images, the PSNR differences caused by using a smaller directional cache are: *Materials* ($32^3$) (28.885dB →28.874dB); *Drums* ($16^3$) (23.745dB →23.836dB); *Lego* ($32^3$) (32.275dB →32.155dB); *Mic* ($32^3$) (31.765dB →31.667dB); *Hotdog* ($32^3$) (34.722dB →34.644dB); *Ficus* ($32^3$) (27.792dB →28.193dB).

### 8.2. Training & Detailed Results

For training, we use the same frequency encoding, noise perturbation and learning rate decay (starting at $5e - 4$) as [27]. For the NeRF 360 Synthetic dataset, we use 64 and 128 samples for the coarse and fine networks, respectively. This becomes 64 and 64 for the LLFF scenes. Note that the coarse networks are discarded and not used in the caching stage - the cached density (or rather, the mesh derived from it) serves as the importance distribution during rendering. We sample 1000 random rays per gradient descent step, and use Adam [12] as our optimizer of choice, with $\beta_1 = 0.9$, $\beta_2 = 0.999$.

In the paper, we compute test metrics on a random subset of 20 images per scene from the test sets of the NeRF 360 Synthetic dataset, and all available test images for the LLFF data. We use the evaluation set only to select the best iteration, which is 300K for all scenes, indicating identical convergence trends as NeRF [27]. For completeness, we show results for all scenes using the *full* number of test images in Table 4. We note that the results and trends are in agreement with the sub-sampled test set. Please see Figure 8 and beyond for a qualitative comparison in addition to the video. We include further ablations on the number of components in Table 5 and Table 6.

### 8.3. Meshing & Rendering

One of the key advantages of working with sparse octrees to store the caches required by our method is that they can serve as a basis for accelerating the rendering process. Using raytracing, we can quickly terminate rays that miss the neurally rendered object(s) all together. For the rays that hit occupied voxels, and so require integrating the volume, we can use raytracing to skip empty space efficiently, saving a significant amount of queries. This matters because caching makes our method memory instead of compute bound. Computing the inner product of the components and weights as well as tracking transmittance for each ray is fast on modern GPUs. Grid look-ups, which require reading from the GPUs RAM on the other hand, are expensive.

While we could use a hierarchical digital differential analyzer such as [28] to determine ray grid intersections, we opt to trace rays against a collision mesh derived from the volume instead, for three reasons. First, this allows us to take advantage of hardware acceleration for the BVH and ray-triangle intersections on modern hardware. Second, we can down/resample or deform the meshes easily using existing tools. Finally, meshes could be used to 'fake' shadows, bouncelight and other effects which can be composited over the neurally rendered content.

We derive the meshes by converting the density cache to a sign distance function (thresholding at $0.0$), optionally downsampling it before meshing with marching cubes, and

Table 4. As in the main text, we compare NeRF to our method when not cached to a grid, and when cached at a high resolution in terms of PSNR, SSIM and LPIPS, and also provide the average speed in ms of our method when cached. The first 8 scenes comprise the dataset of [27] at $800^2$ pixels, the last 8 scenes the LLFF dataset [26] at $504 \times 378$ pixels. We use 8 components and a $1024^3$ cache for the NeRF Synthetic 360 scenes (except for the ship scene where we use $768^3$), and 6 components at $768^3$ for LLFF.

| Scene | NeRF | | | Ours - No Cache | | | Ours - Cache | | | Speed |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | PSNR↑ | SSIM↑ | LPIPS↓ | |
| **Drums** | 24.58dB | 0.92 | 0.08 | 23.868dB | 0.91 | 0.084 | 23.745dB | 0.913 | 0.083 | 4.5ms |
| **Ship** | 27.21dB | 0.83 | 0.17 | 27.241dB | 0.839 | 0.155 | 27.685dB | 0.805 | 0.192 | 4.9ms |
| **Mic** | 30.85dB | 0.97 | 0.03 | 30.274dB | 0.973 | 0.023 | 31.765dB | 0.977 | 0.022 | 2.6ms |
| **Ficus** | 28.86dB | 0.96 | 0.03 | 27.037dB | 0.948 | 0.034 | 27.792dB | 0.954 | 0.031 | 3.7ms |
| **Chair** | 31.16dB | 0.96 | 0.04 | 30.967dB | 0.96 | 0.032 | 32.322dB | 0.966 | 0.032 | 2.6ms |
| **Lego** | 30.41dB | 0.95 | 0.03 | 31.076dB | 0.955 | 0.023 | 32.275dB | 0.964 | 0.022 | 5.5ms |
| **Hotdog** | 34.7dB | 0.97 | 0.03 | 34.21dB | 0.97 | 0.029 | 34.722dB | 0.973 | 0.031 | 4.9ms |
| **Materials** | 28.93dB | 0.94 | 0.03 | 28.567dB | 0.943 | 0.034 | 28.885dB | 0.947 | 0.034 | 3.9ms |
| **Fern*** | 26.84dB | 0.86 | 0.09 | 26.325dB | 0.853 | 0.105 | 25.006dB | 0.822 | 0.131 | 1.7ms |
| **Flower*** | 28.32dB | 0.89 | 0.06 | 28.916dB | 0.912 | 0.043 | 27.012dB | 0.878 | 0.059 | 1.3ms |
| **Fortress*** | 32.7dB | 0.93 | 0.03 | 32.946dB | 0.937 | 0.021 | 31.256dB | 0.913 | 0.043 | 0.8ms |
| **Horns*** | 28.78dB | 0.9 | 0.07 | 29.748dB | 0.925 | 0.044 | 26.847dB | 0.889 | 0.07 | 1.2ms |
| **Leaves*** | 22.43dB | 0.82 | 0.1 | 22.53dB | 0.833 | 0.095 | 21.3dB | 0.787 | 0.122 | 1.5ms |
| **Orchids*** | 21.36dB | 0.75 | 0.12 | 21.204dB | 0.747 | 0.126 | 20.356dB | 0.712 | 0.137 | 1.6ms |
| **Room*** | 33.59dB | 0.96 | 0.04 | 33.702dB | 0.961 | 0.034 | 30.301dB | 0.942 | 0.057 | 1.6ms |
| **TRex*** | 27.74dB | 0.92 | 0.05 | 28.292dB | 0.933 | 0.04 | 26.204dB | 0.903 | 0.063 | 1.4ms |

Table 5. Influence of the number of components and grid resolution on PSNR and memory required for caching the *Horns* ('Triceratops') scene. Note how more factors also increase grid sparsity in this case. 6 components are a reasonable amount for LLFF data.

| Factors | No Cache | | $256^3$ | | $384^3$ | | $512^3$ | | $768^3$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR↑ | Memory | PSNR↑ | Memory | PSNR↑ | Memory | PSNR↑ | Memory | PSNR↑ | Memory |
| **4** | 29.56dB | - | 22.94dB | 0.39Gb | 24.64dB | 0.8Gb | 25.62dB | 1.56Gb | 26.52dB | 4.18Gb |
| **6** | 29.75dB | - | 23.08dB | 0.56Gb | 24.86dB | 1.14Gb | 25.89dB | 2.22Gb | 26.85dB | 6.11Gb |
| **8** | 29.82dB | - | 23.02dB | 0.72Gb | 24.79dB | 1.45Gb | 25.81dB | 2.81Gb | 26.76dB | 7.83Gb |

optionally simplifying the resulting geometry using standard techniques. We show the resulting geometry for the *Lego* scene in Figure 7. Note that more aggressive thresholding of the volume or remeshing is easily possible and can lead to significant performance benefits in practise as more rays can be culled early and mesh complexity reduced.

For rendering, we follow the same method as NeRF, using the Beer-Lambert Law [7] to model radiance extinction. Once the contribution of new samples for a ray reaches 0.001, we terminate a ray's rendering kernel. Note that this can be relaxed for greater performance. Another way to accelerate the rendering process is to visit fewer voxels based on radiance or transmittance. For all paper experiments however, we never skip occupied space, or optimise any parameters per scene.

## 8.4. Cache size calculation

In this section we show the formulas we used to estimate cache sizes $M_{NeRF}, M_{FastNeRF}$ for NeRF and FastNeRF respectively.

$$M_{NeRF} = \alpha(s_\sigma + s_{rgb})k^3 l^2, \quad (6)$$

where $s_\sigma, s_{rgb}$ are the sizes of the stored transparency and RGB values in bits and $\alpha \in \langle 0, 1 \rangle$ is the inverse volume sparsity, where $\alpha = 1$ would indicate a dense volume.

$$M_{FastNeRF} = \alpha \left( (D \cdot s_{abc} + s_\sigma)k^3 \right) + D \cdot s_\beta l^2, \quad (7)$$

where $s_{uvw}, s_\beta$ are the sizes of the stored $(u_i, v_i, w_i)$ values and the weights $\beta_i$.

Table 6. Influence of the number of components and grid resolution on PSNR and memory required for caching the *Fortress* ('Minas Tirith') scene. Note how more factors also increase grid sparsity in this case. 6 components are a reasonable amount for LLFF data.

| Factors | No Cache | | $256^3$ | | $384^3$ | | $512^3$ | | $768^3$ | |
|---|---|---|---|---|---|---|---|---|---|---|
| | *PSNR↑* | *Memory* | *PSNR↑* | *Memory* | *PSNR↑* | *Memory* | *PSNR↑* | *Memory* | *PSNR↑* | *Memory* |
| **4** | 32.81dB | - | 26.54dB | 0.42Gb | 28.61dB | 0.93Gb | 29.83dB | 1.88Gb | 30.99dB | 5.43Gb |
| **6** | 32.95dB | - | 26.63dB | 0.6Gb | 28.75dB | 1.31Gb | 30.01dB | 2.66Gb | 31.26dB | 7.83Gb |
| **8** | 33.03dB | - | 26.66dB | 0.79Gb | 28.81dB | 1.72Gb | 30.1dB | 3.47Gb | 31.36dB | 10.29Gb |



Figure 7. Collision Meshes for the *Lego* scene used for our raytraced rendering algorithm. From left to right, top to bottom, these are derived from the following grid resolutions $256^3$, $384^3$, $512^3$, $768^3$, $1024^3$. Note that we threshold the density at 0, but being less conservative can give meshes good enough to compute shadows or bounce light. This is seen in the bottom-right image, where the identical volume is thresholded at 10.0 instead (best viewed zoomed in).

For $k = l = 1024$, $(r, g, b)$ stored as individual bytes and all other values stored as half-precision floats the cache size would be

$$M_{NeRF} = \alpha \cdot 5,629,499,534,213,120 \approx \alpha \cdot 5.6PB, \quad (8)$$

for NeRF and

$$M_{factor} = \alpha \cdot 105,226,698,752 + 33,554,532 \approx \alpha \cdot 105GB, \quad (9)$$

for FastNeRF.

Even for highly sparse volumes, where $\alpha = 0.01$, the cache would take nearly 60TB, which exceeds the mem-ory capacity of consumer-grade hardware. The high memory requirement makes this approach impractical when used with the standard NeRF architecture.

## 8.5. Cache size calculation

In this section we show the formulas we used to estimate cache sizes $M_{NeRF}$, $M_{FastNeRF}$ for NeRF and FastNeRF respectively.

$$M_{NeRF} = \alpha(s_\sigma + s_{rgb})k^3l^2, \quad (10)$$

where $s_\sigma, s_{rgb}$ are the sizes of the stored transparency and RGB values in bits and $\alpha \in \langle 0, 1 \rangle$ is the inverse volume

13

sparsity, where $\alpha = 1$ would indicate a dense volume. Just as in the main paper $k, l$ correspond to the number of bins per dimension in the position and direction dependent caches respectively.

$$M_{FastNeRF} = \alpha \left( (D \cdot s_{uvw} + s_\sigma)k^3 \right) + D \cdot s_\beta l^2, \quad (11)$$

where $s_{uvw}, s_\beta$ are the sizes of the stored $(u_i, v_i, w_i)$ values and the weights $\beta_i$.

For $k = l = 1024, D = 8, (r, g, b)$ stored as individual bytes and all other values stored as half-precision floats the cache size would be

$$M_{NeRF} = \alpha \cdot 5,629,499,534,213,120 \approx \alpha \cdot 5.6PB, \quad (12)$$

for NeRF and

$$M_{factor} = \alpha \cdot 53,687,091,200 + 16,777,216 \approx \alpha \cdot 54GB, \quad (13)$$

for FastNeRF.

### 8.6. Details on Applications

In our proof-of-concept telepresence scenario we use a deformation field network [32] $\mathcal{F}_{deform}$ that modifies the input sample positions. This network is similar to the position-dependent network used in FastNeRF, but smaller - a 6-layer MLP with 64 units in each layer. The input is a sample position and an expression vector. The sample position is processed with positional encoding identically to how the FastNeRF input position is processed. The output is an offset that is applied to the input sample position. When training with $\mathcal{F}_{deform}$ we add an L2 regularizer that constrains $\mathcal{F}_{deform}$ to be an identity transform if a neutral expression is passed as input. We find that this solution stabilizes training and improves results.

In addition to $\mathcal{F}_{deform}$, we also use a non-trainable deformation field $\mathcal{F}_{bone}$ that moves the samples in the head region according to the movement of the 3D model bones that represent the shoulders and the neck. This additional deformation accounts for large movements of the head and the shoulders, which reduces the load on $\mathcal{F}_{deform}$ and improves quality. The two deformation fields are composed with the position-dependent network $\mathcal{F}_{pos}$ as follows:

$$\{\sigma, (\boldsymbol{u}, \boldsymbol{v}, \boldsymbol{w})\} = \mathcal{F}_{pos}\left(\mathcal{F}_{deform}(\mathcal{F}_{bone}(\boldsymbol{p}, \boldsymbol{e}), \boldsymbol{e})\right), \quad (14)$$

where $\boldsymbol{e}$ is the expression vector. While we could use $\mathcal{F}_{bone}$ at test-time as well, we remove it to improve performance.

The training procedure in this scenario is identical to that used in FastNeRF but with 64 samples in the coarse stage and an additional 64 in the fine stage. The position-dependent network has 384 units in each layer, while the view-dependent network uses only 32 units and 2 layers. The view-dependent network is very small as we believe the scene's illumination is simple and fairly uniform.

The main limitation of this approach in terms of speed is the need to call $\mathcal{F}_{deform}$ for all the input samples. In order to mitigate this, at test-time we implement a grid-based sample pruner. For each position in a grid around the face region we compute the density $\sigma$

$$\sigma = \mathcal{F}_{pos}\left(\mathcal{F}_{deform}(\boldsymbol{p}, \boldsymbol{e})\right) \quad (15)$$

for 50 randomly chosen samples of $\boldsymbol{e}$. For each grid position we record the minimum and maximum encountered $\sigma$ in a sparse volume. At test-time we use this volume to prune the inputs to $\mathcal{F}_{deform}$ where the maximum $\sigma$ is 0 and where the rays would get saturated, which we estimate using the minimum density values from the volume. This approach reduces the runtime of $\mathcal{F}_{deform}$ by more than a factor of 2.

To further reduce the impact of $\mathcal{F}_{deform}$ on the runtime we only use 43 samples in the coarse stage and remove the fine stage altogether at test time. To offset this low sample count we add additional samples that are linearly interpolated between the samples generated by the deformation network

$$\boldsymbol{p}_{deform} = \mathcal{F}_{deform}(\boldsymbol{p}, \boldsymbol{e}). \quad (16)$$

The additional samples can be evaluated very cheaply as they can be looked-up in the FastNeRF cache. Note that since the deformation network changes the positions of individual samples along the rays, the rays are no longer straight. This means that in this scenario we cannot use the hardware-accelerated ray tracing procedure described in the main paper, though the pruning method described above serves a similar role.

Even though the steps above significantly reduce the runtime of $\mathcal{F}_{deform}$, we are still limited to rendering $300 \times 300$ pixel images with a reduced sample count if we want to maintain 30FPS with the deformation network. We observed that if this framerate constraint was to be disregarded, the approach described above is able to generate images at significantly better quality. Thus, we believe that a faster deformation approach remains an important goal for future work.
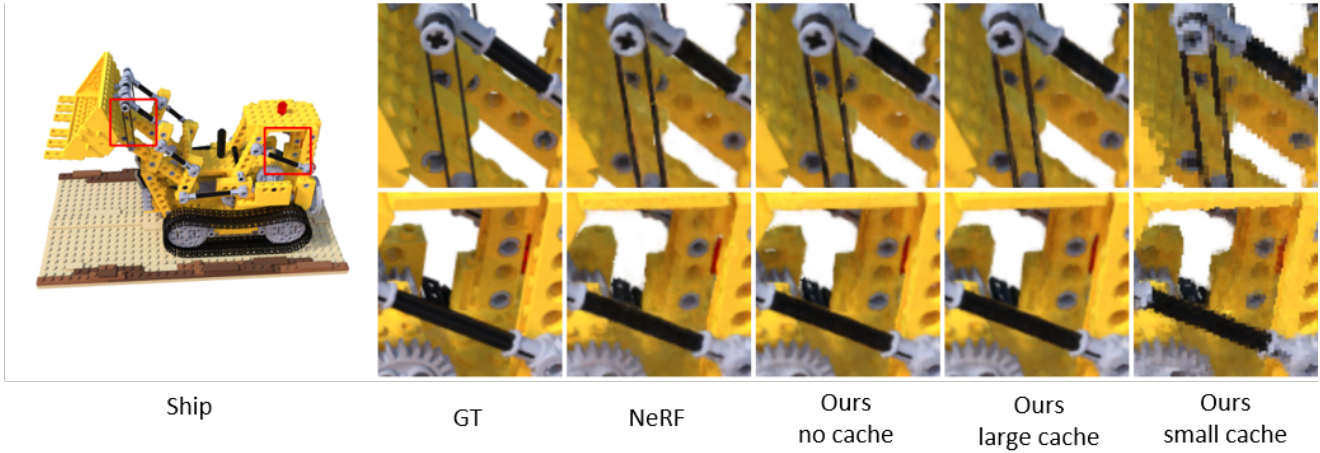
Figure 8. Qualitative comparison of our method vs NeRF on the *Lego* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.
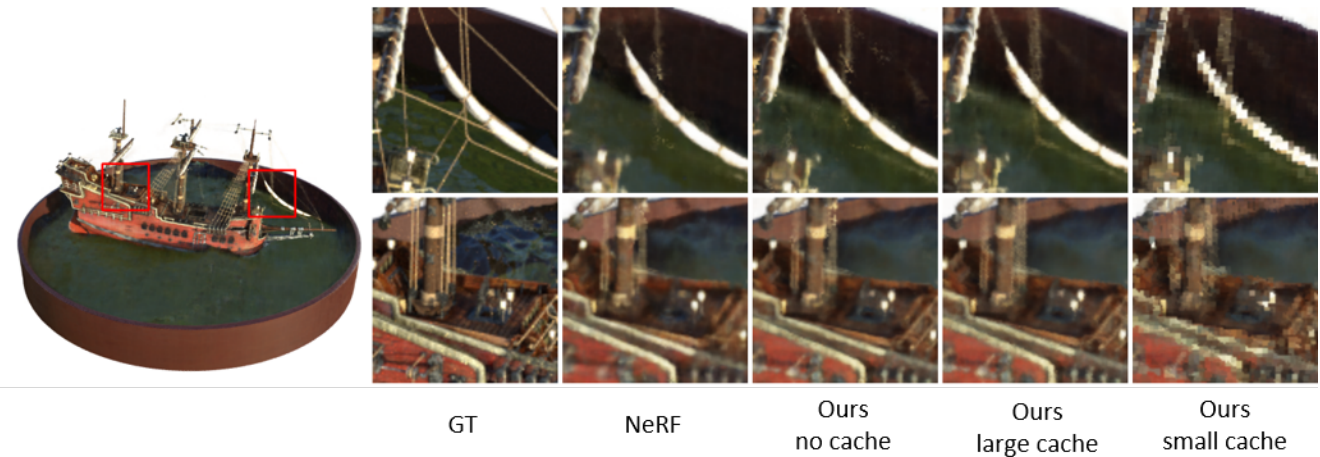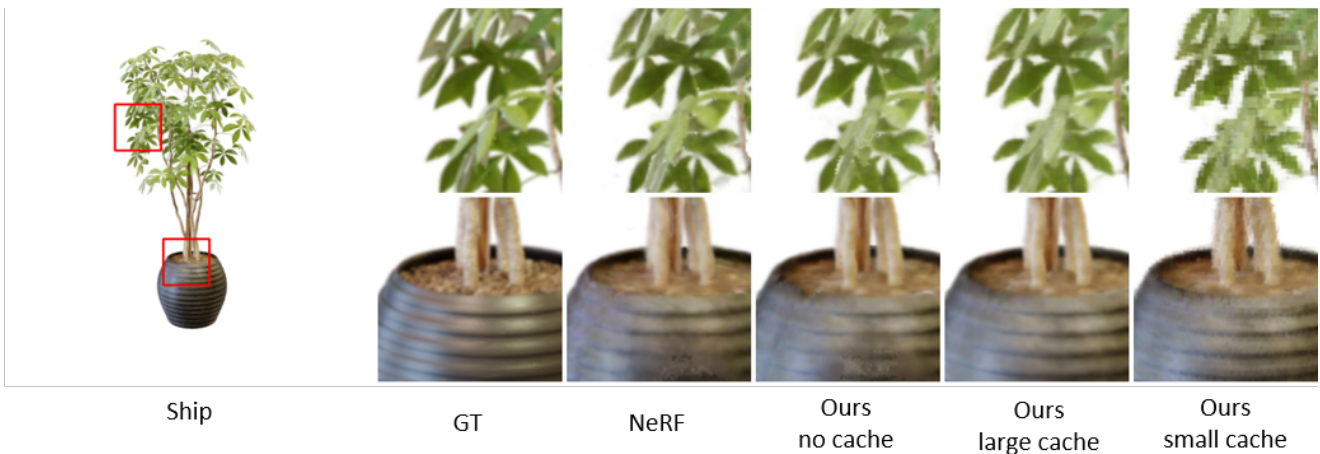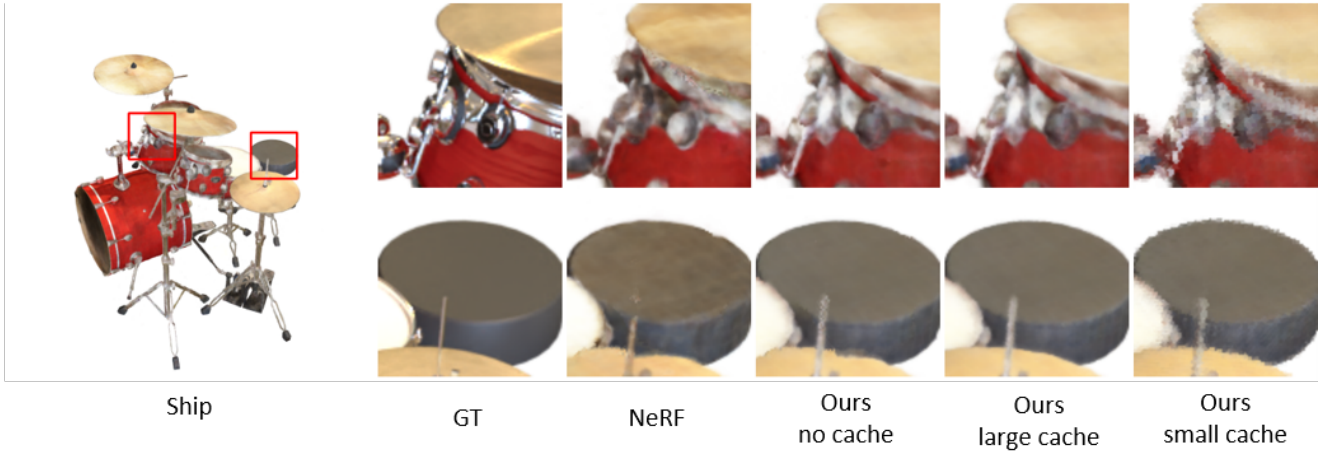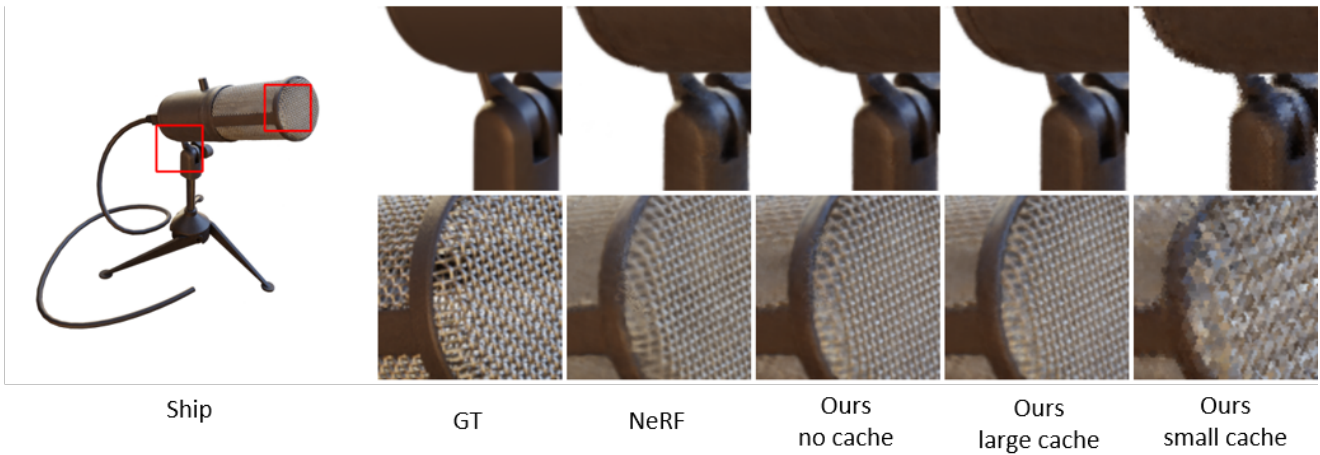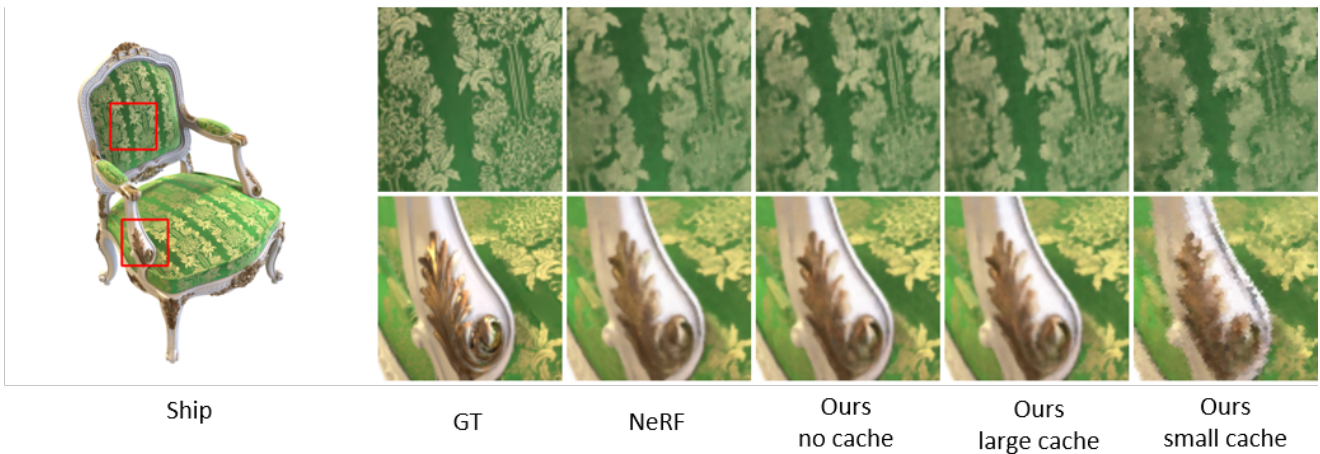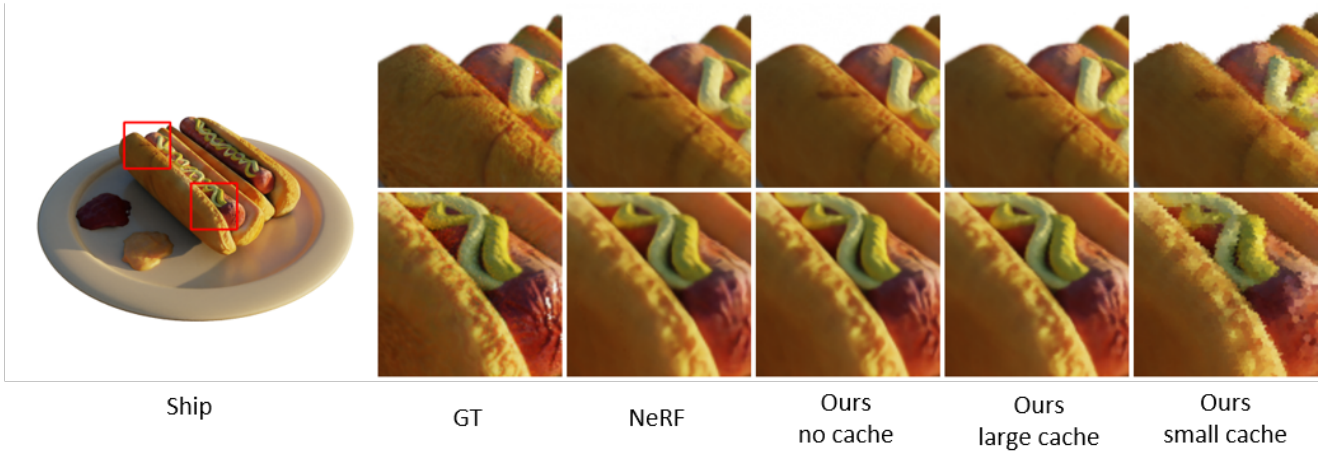


Figure 9. Qualitative comparison of our method vs NeRF on the *Ship* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.
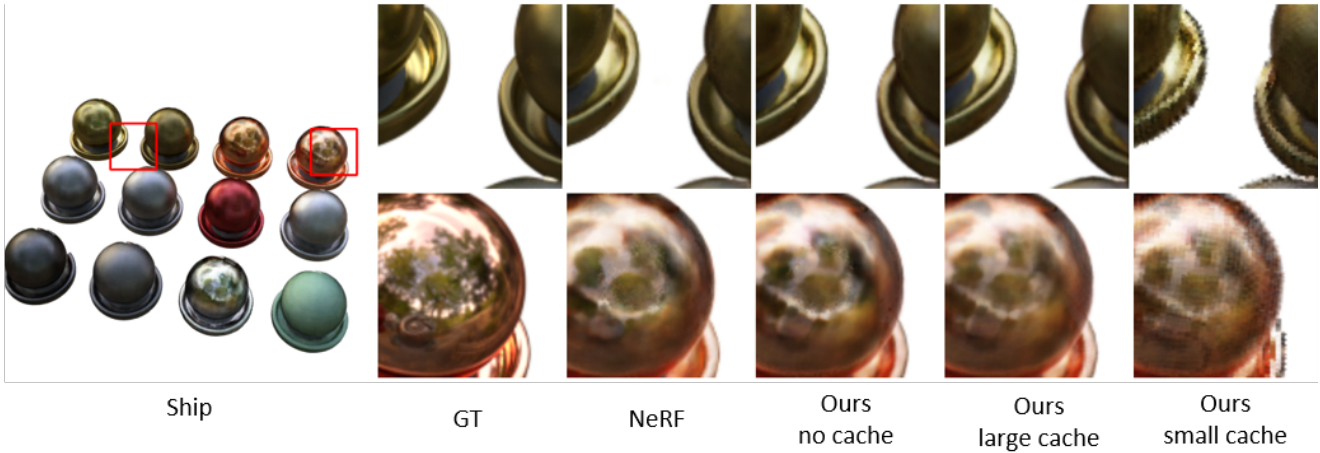


Figure 10. Qualitative comparison of our method vs NeRF on the *Ficus* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.

15

Figure 11. Qualitative comparison of our method vs NeRF on the *Drums* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.
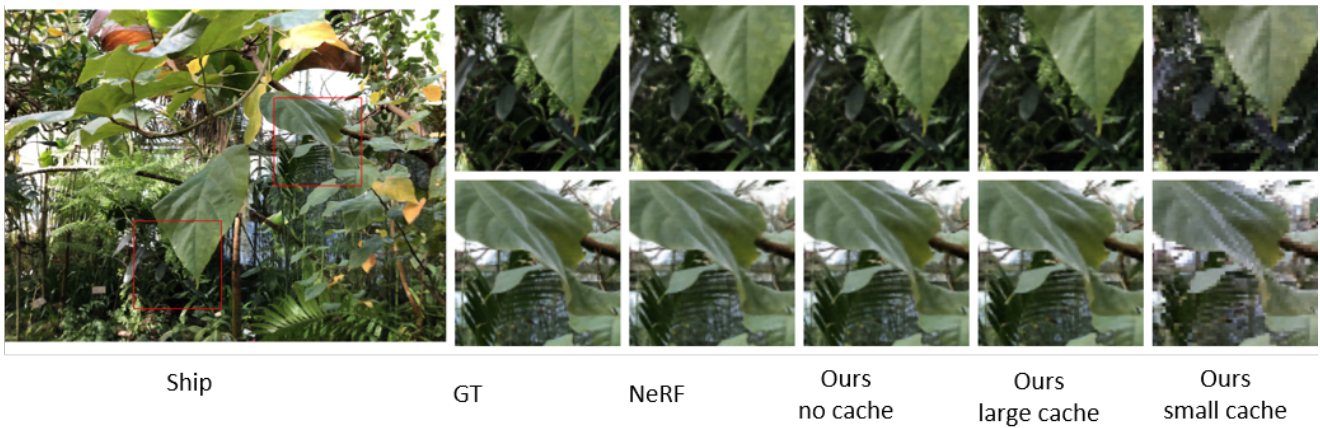


Figure 12. Qualitative comparison of our method vs NeRF on the *Mic* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.
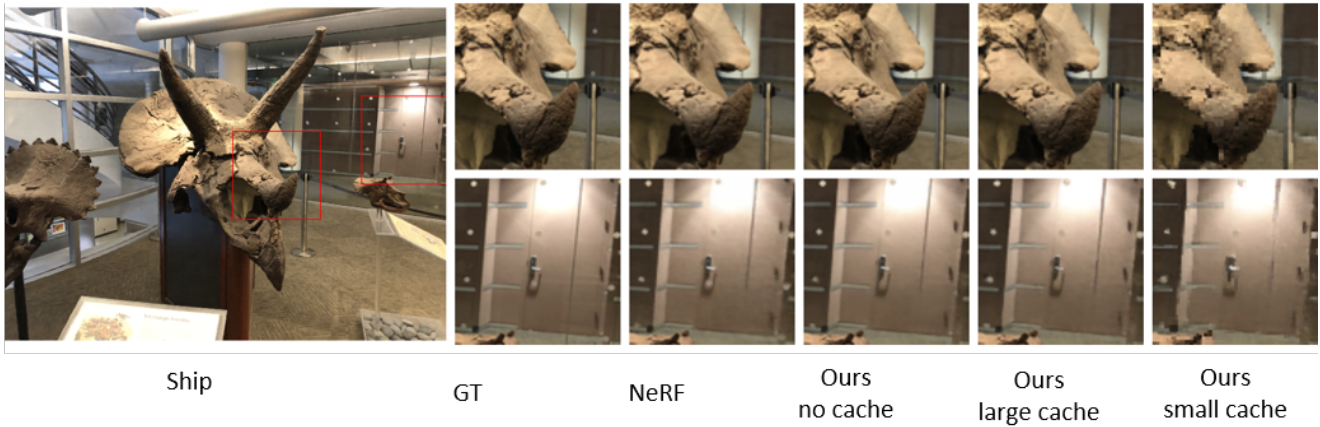


Figure 13. Qualitative comparison of our method vs NeRF on the *Chair* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.

16

Figure 14. Qualitative comparison of our method vs NeRF on the *Hotdog* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.



Figure 15. Qualitative comparison of our method vs NeRF on the *Materials* scene from the dataset of [27] at $800^2$ pixels using 8 components. *Small cache* refers to our method cached at $256^3$, and *large cache* at $1024^3$.



Figure 16. Qualitative comparison of our method vs NeRF on the *Leaves* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.
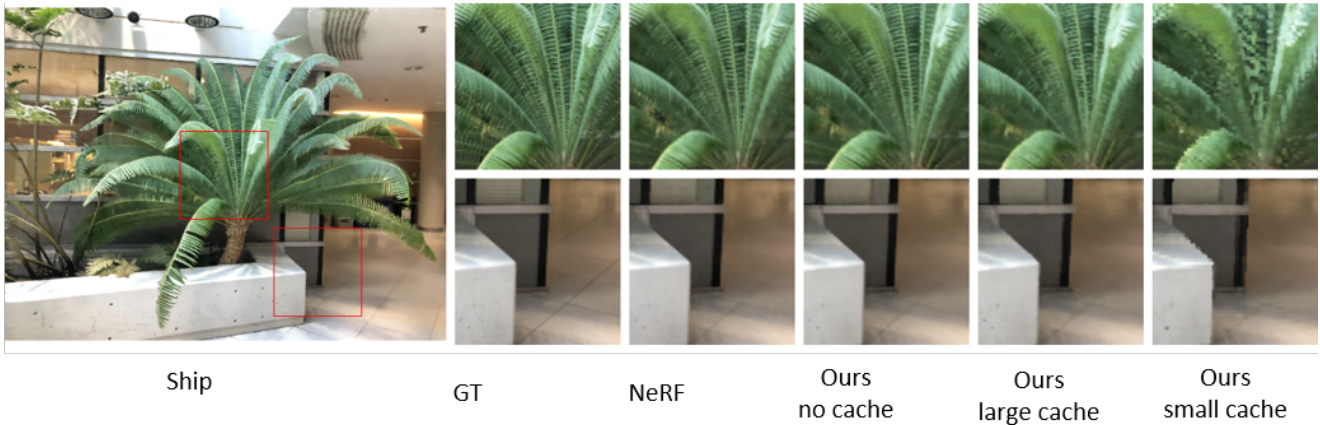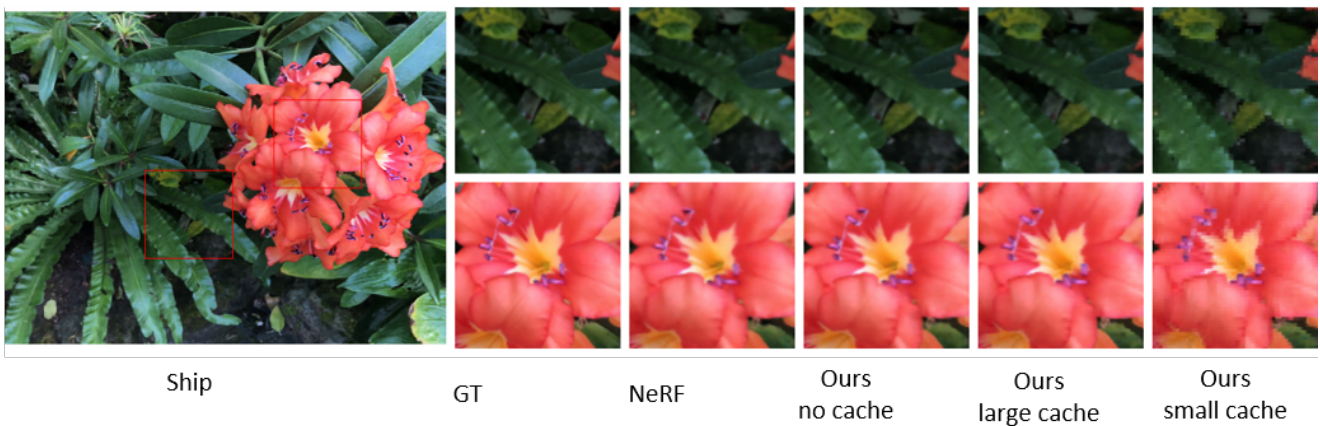
17

Figure 17. Qualitative comparison of our method vs NeRF on the *Horns* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.



Figure 18. Qualitative comparison of our method vs NeRF on the *Fern* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.



Figure 19. Qualitative comparison of our method vs NeRF on the *Flower* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.
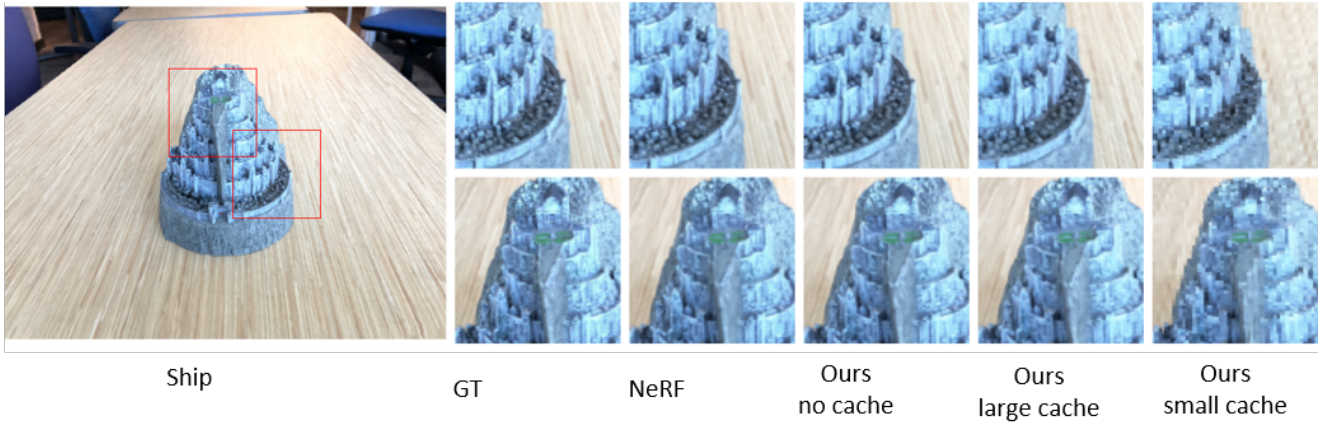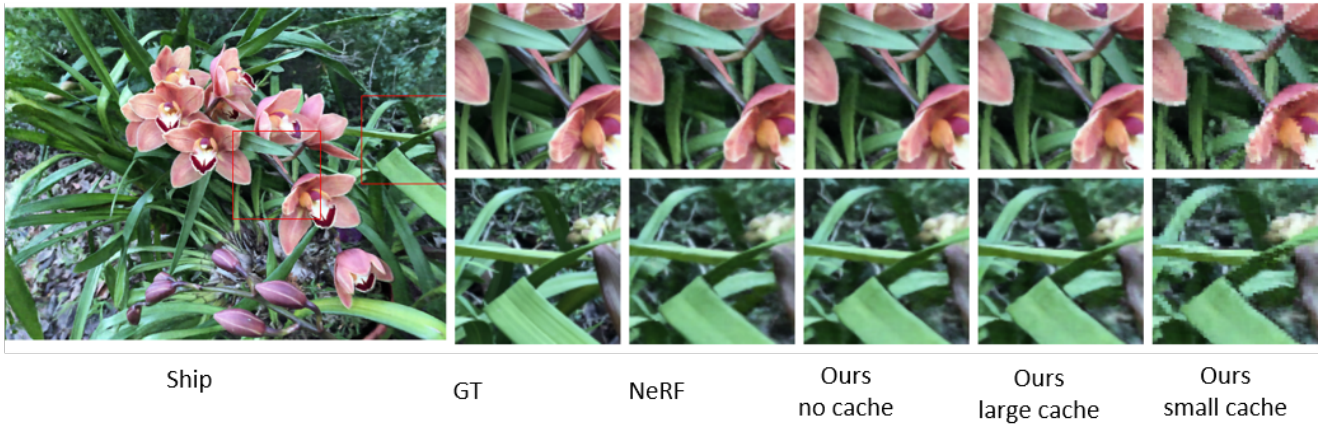
Figure 20. Qualitative comparison of our method vs NeRF on the *Fortress* ('Minas Tirith') scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.



Figure 21. Qualitative comparison of our method vs NeRF on the *Orchids* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.
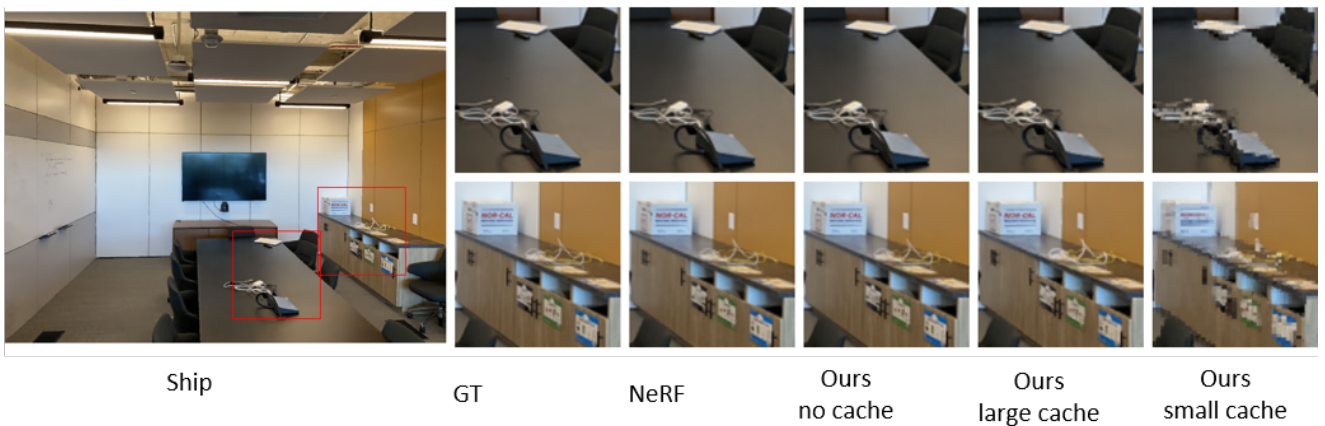


Figure 22. Qualitative comparison of our method vs NeRF on the *Room* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.
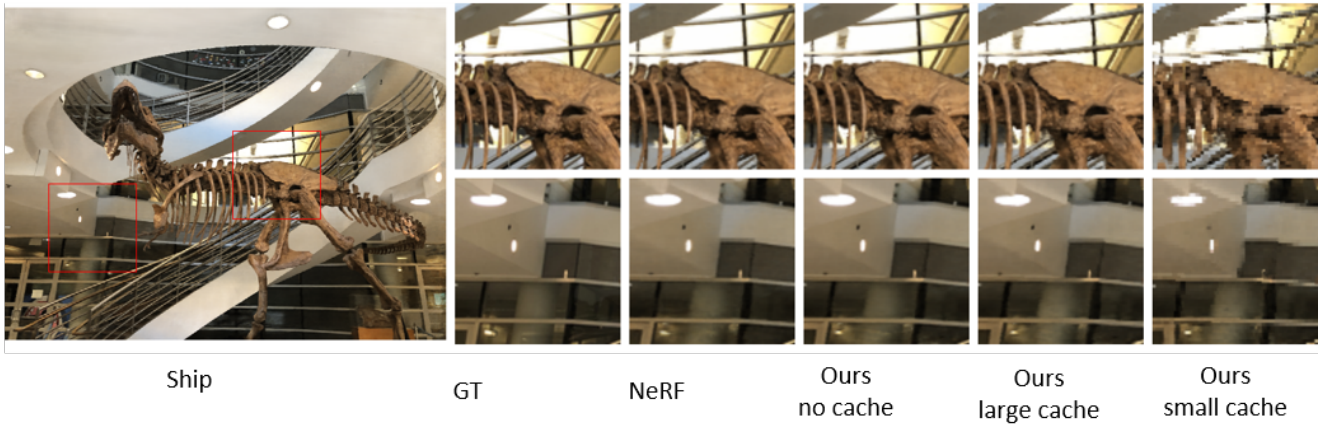
Figure 23. Qualitative comparison of our method vs NeRF on the *TRex* scene from the dataset of [26] at $504 \times 378$ pixels using 6 factors. *Small cache* refers to our method cached at $256^3$, and *large cache* at $768^3$.