

InterNeRF: Scaling Radiance Fields via Parameter Interpolation

Clinton Wang^{*1} Peter Hedman² Polina Golland¹ Jonathan T. Barron² Daniel Duckworth³

¹MIT CSAIL ²Google Research ³Google DeepMind

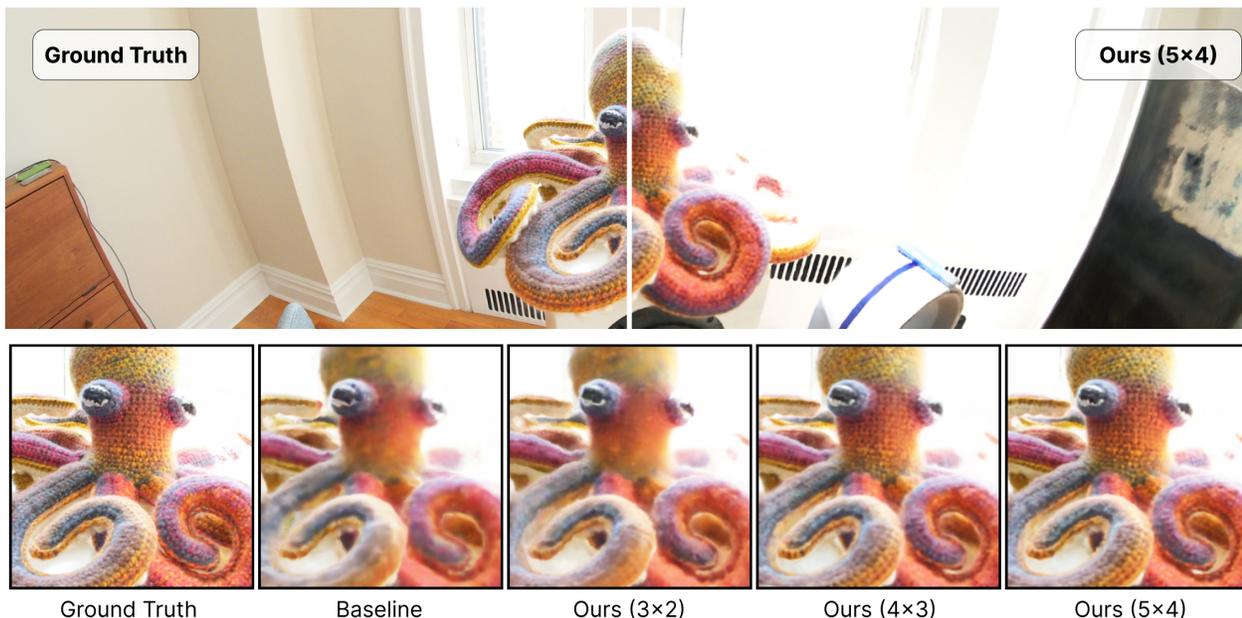


Figure 1. InterNeRF achieves state-of-the-art reconstruction quality on large, multi-room scenes via parameter interpolation. Parameters are anchored to locations in space and interpolated based on camera position. As a result, InterNeRF achieves high geometric and texture accuracy at centimeter scale.

Abstract

Neural Radiance Fields (NeRFs) have unmatched fidelity on large, real-world scenes. A common approach for scaling NeRFs is to partition the scene into regions, each of which is assigned its own parameters. When implemented naively, such an approach is limited by poor test-time scaling and inconsistent appearance and geometry. We instead propose InterNeRF, a novel architecture for rendering a target view using a subset of the model’s parameters. Our approach enables out-of-core training and rendering, increasing total model capacity with only a modest increase to training time. We demonstrate significant improvements in multi-room scenes while remaining competitive on standard benchmarks.

^{*}Work done as a student researcher at Google.

1. Introduction

Neural Radiance Fields (NeRFs) are a class of powerful, high-fidelity representations for 3D reconstruction achieving unparalleled reconstruction accuracy. Unlike many approaches, NeRF models are straightforward to train and resilient to local minima while delivering impressive quality across a wide variety of scenes. Nevertheless, reconstruction quality is inherently limited by model capacity, and new methods are needed to scale beyond the memory and compute limitations of present-day hardware.

A natural way to increase capacity is to spatially partition network parameters based on geometry [9, 11, 15, 19] or camera [14] location. Geometry-based partitioning divides the scene into multiple regions, which each have their own set of parameters. While a fraction of model parameters are required to render a target ray, the number and choice of pa-

parameters varies depending on ray origin and direction and occlusions within the scene itself. Camera-based partitioning, on the other hand, assigns parameters to possible query cameras. While rendering is typically simpler and more efficient, inconsistencies arise when multiple parameter sets redundantly represent the same scene content.

We introduce InterNeRF, a high-capacity NeRF architecture tailored to large, multi-room scenes. Key to our approach is the concept of *camera-centric parameter interpolation*: the efficient interpolation of network parameters based on camera origin. In particular, we designate a subset of model parameters as spatially-partitioned, loading and unloading parameters relevant to the active camera region during training. At rendering time, the appropriate model parameters are loaded based on camera origin. Given sufficient training time, our method outperforms a state-of-the-art baseline by a wide margin on the Zip-NeRF dataset.

2. Preliminaries

NeRF: A NeRF [7] is a neural network that represents a 3D scene as a continuous mapping from a 3D spatial location \mathbf{x}_i to a volumetric density τ_i and color \mathbf{c}_i . Novel images of the scene are rendered by tracing rays through this neural representation, where each ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is rendered according to the Gaussian quadrature approximation to the volume rendering equation,

$$\mathbf{C} = \sum_i T_i (1 - \exp(-\tau_i \delta_i)) \mathbf{c}_i \quad (1)$$

$$T_i = \exp\left(\sum_{j=1}^{i-1} -\tau_j \delta_j\right) \text{ where } \delta_i = t_{i+1} - t_i. \quad (2)$$

Ray intervals are partitioned into non-overlapping intervals $\{[t_i, t_{i+1}]\}$, and a multilayer perceptron (MLP) is trained to estimate τ_i and \mathbf{c}_i of these intervals from \mathbf{x}_i (taken to be the center of each interval) [7]. This process is repeated, where a *proposal network* learns to map a uniform sampling of distances $\{t_i\}$ to a “coarse” set of densities, which are then resampled to produce new distances that are concentrated around scene content [3].

Instant-NGP and Zip-NeRF: Instant-NGP [8] introduced a learned, multi-resolution datastructure for efficiently featureizing spatial coordinates throughout the scene. At coarse scales in the datastructure, spatial coordinate \mathbf{x}_i is trilinearly interpolated into a dense, multi-channel 3D grid to produce feature vectors (which we refer to as *grid features*), and at fine scales \mathbf{x}_i is trilinearly interpolated into a 3D grid backed by a hash table to produce a feature vector (which we call *hash features*). Instant-NGP works by querying a spatial coordinate across all scales and concatenating the interpolated outputs to construct a feature vector. We write this interpolation and concatenation as: $\mathbf{z}_i = \text{NGP}(\mathbf{x}_i)$.

A small MLP called the “geometry MLP” then takes these feature vectors and predicts a scalar density value τ_i , and a larger MLP takes both \mathbf{z}_i and the viewing direction of the ray \mathbf{d} as input to predict radiance $\mathbf{c}_i = \text{MLP}(\mathbf{z}_i, \mathbf{d})$.

Building on mip-NeRF [2] and Instant-NGP, Zip-NeRF [4] uses multisampling and reweighting to parameterize the conical sub-frusta along each ray using NGPs. This strategy achieves the anti-aliasing and scale-reasoning provided by mip-NeRF while being significantly faster to optimize and evaluate. Zip-NeRF still attains the highest quality in novel view synthesis, outperforming 3D Gaussian Splatting [6] on photos acquired with rectilinear lenses and additionally able to accommodate distorted lenses (which 3DGS cannot).

3. InterNeRF

We propose Interpolated NeRF (InterNeRF), a scalable, out-of-core approach for increasing model capacity without a corresponding increase in memory usage. We adopt the concept of *spatially-partitioned* model parameters: each parameter is broadly categorized as either INTERPOLATED or SHARED, with INTERPOLATED parameters varying with camera origin and SHARED remaining unchanged. To render a particular camera, we interpolate INTERPOLATED parameters within a local neighborhood and perform a forward pass similar to typical NeRF model.

While interpolated parameters provide the model the freedom to specialize to relevant camera views, they also permit origin-dependent geometry, and geometric coherence can suffer when all parameters are interpolated. Thus the role of the SHARED parameters is to maintain a stable coarse geometry, and we choose to share proposal network parameters and all NGP grid features.

Parameter grid: To determine the anchor locations for parameter sets, we begin by establishing a 2D axis-aligned bounding box containing all training cameras. The bounding box is partitioned into $N_x \times N_y$ non-isotropic grid cells. We assume that scenes are largely planar and thus omit a partitioning along the vertical axis. Cameras are then assigned to cells based on their origin (see Fig. 2). Note that only cells with a sufficient number of training cameras are instantiated; if a camera’s origin lies outside of an instantiated cell, cameras are reassigned to the nearest active cell instead.

Parameter mixing: Similar to Block-NeRF, we assign each vertex in the parameter grid its own *parameter set* for interpolation. However, Block-NeRF applies nearest neighbor interpolation, which introduces discontinuities at the boundaries between parameter sets and necessitates post-processing heuristics like visibility maps and image-space interpolation to avoid “popping” artifacts. We instead explore *bilinear* interpolation, which naturally results in a continuous field of model parameters.

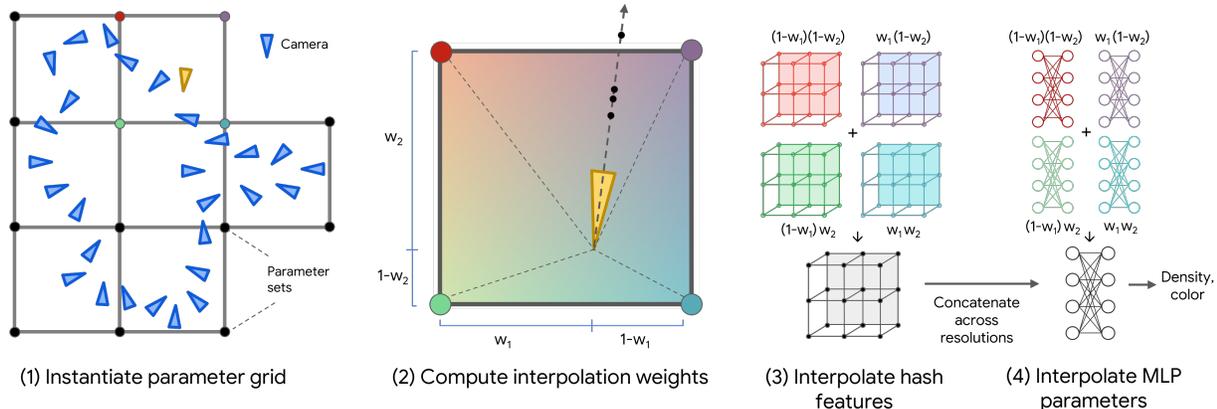


Figure 2. **Our framework.** (1) We partition the scene into a parameter grid and assign training cameras to each cell based on its origin. (2) For a given (training or target) view, we obtain mixing weights as the bilinear interpolation coefficients based on camera origin. (3) Each query point along the ray is used to index into a multi-resolution set of grid features per parameter set, with either explicit assignment or a hash table. The mixing weights are applied here to yield a single set of features. (4) Each parameter set also has its own set of MLP weights, which are combined using the same mixing weights to form a new MLP.

We implement bilinear interpolation layer-by-layer, performing a forward pass of each layer with four parameter sets followed by interpolation of their outputs. In the NGP component, four sets of hash features are retrieved and interpolated at each resolution; in the MLP, each fully connected layer performs four multiply-adds followed by interpolation. This approach avoids constructing per-sample model parameters, which would significantly increase memory usage.

Cell-by-cell training: To decouple memory usage from parameter count, we optimize four parameter sets at a time corresponding to a single grid cell. We optimize cells serially in a round robin fashion, loading and unloading partitioned model parameters as needed. When optimizing a cell, we by default use training cameras that lie within the same cell. Rather than optimizing cells for a fixed number of iterations, we vary the number of iterations proportionally to the number of cameras assigned to it. We find this improves quality and reduces floating artifacts.

Ray reassignment: To increase the amount of training signal each cell receives, we incorporate cameras from other cells during optimization. In particular, each training batch is constructed such that a fixed percentage p_r of camera rays are sourced from neighbors at most K_r cells away. The interpolation weights are determined by projecting the camera origin to the closest point in the cell. We find that this strategy strongly reduces floating geometry immediately outside of the training camera frustum.

4. Experiments

Datasets: We evaluate our approach on two datasets: four large multi-room scenes introduced by Zip-NeRF [4] and

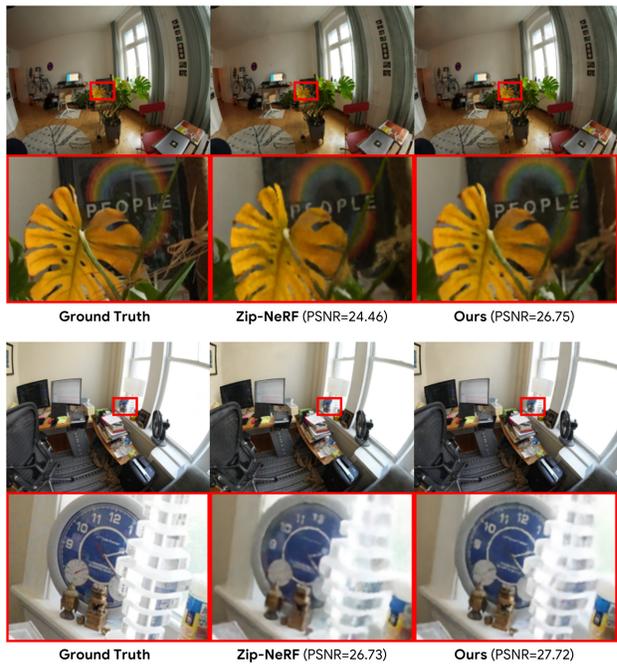


Figure 3. Images from BERLIN and NYC in the Zip-NeRF dataset, rendered by Zip-NeRF and InterNeRF with a 5×4 grid.

the four indoor and five outdoor scenes of the mip-NeRF 360 dataset [3]. All models were trained from scratch using the original, distorted photos with every 8th image held out for test. To illustrate the higher fidelity of our approach, we use images at double their typical resolution: 1752×1168 for Zip-NeRF and 2456×1632 or 3120×2080 for mip-NeRF 360. All camera parameters have been estimated with COLMAP at full resolution [12]. We note that this dataset

Method / Steps	PSNR				SSIM				LPIPS				Wall Time @ 400K
	50K	100K	200K	400K	50K	100K	200K	400K	50K	100K	200K	400K	
Zip-NeRF	25.71	25.82	25.84	25.83	0.796	0.799	0.800	0.799	0.373	0.368	0.365	0.366	6.8h
Ours (3×2)	25.83	25.98	26.07	26.14	0.801	0.803	0.805	0.805	0.372	0.367	0.363	0.362	17.5h
Ours (4×3)	25.95	26.24	26.35	26.31	0.809	0.815	0.818	0.818	0.358	0.346	0.341	0.339	19.9h
Ours (5×4)	25.61	26.28	26.42	26.65	0.806	0.818	0.823	0.826	0.358	0.341	0.331	0.326	23.4h

Table 1. Quantitative evaluation on the Zip-NeRF dataset. Quality increases monotonically with scene grid resolution and training steps for our method while Zip-NeRF saturates. To maintain comparable memory usage, we provide Zip-NeRF with 4× more hash grid parameters than our interpolated models.

Method / Steps	PSNR				SSIM				LPIPS				Wall Time @ 200K
	25K	50K	100K	200K	25K	50K	100K	200K	25K	50K	100K	200K	
Zip-NeRF	28.15	28.28	28.34	28.34	0.808	0.813	0.816	0.816	0.313	0.305	0.299	0.296	2.7h
Ours (2×2)	28.23	28.35	28.42	28.42	0.801	0.808	0.812	0.815	0.327	0.315	0.308	0.304	5.2h
Ours (3×3)	28.08	28.33	28.43	28.49	0.795	0.807	0.814	0.817	0.340	0.322	0.312	0.306	7.1h
Ours (4×4)	27.86	28.21	28.38	28.46	0.779	0.797	0.809	0.814	0.363	0.338	0.322	0.312	10.7h

Table 2. Quantitative evaluation on four scenes from the mip-NeRF 360 dataset: BICYCLE, GARDEN, COUNTER, and BONSAI. PSNR and SSIM are competitive with Zip-NeRF as the number of training steps increases. Here we again provide Zip-NeRF with 4× more hash grid parameters than our interpolated models.

	PSNR	SSIM	LPIPS	Wall Time
Zip-NeRF	27.55	0.794	0.292	2.73 Hr
Ours (2x2)	27.66	0.795	0.299	5.21 Hr
Ours (3x3)	27.69	0.796	0.300	7.12 Hr
Ours (4x4)	27.64	0.793	0.307	10.68 Hr

Table 3. Quantitative evaluation on the full mip-NeRF 360 dataset. Models are trained for 200K steps.

is not amenable rasterization methods such as 3D Gaussian Splatting, which assume a pinhole camera model.

Baseline: We compare our method to Zip-NeRF, a state-of-the-art NeRF method for reconstruction of large indoor spaces. Our implementation uses three networks: two for proposing ray intervals and one for color and density prediction. All three networks consist of a multiresolution hash grid with up to 2^{21} entries followed by a per-interval geometry MLP with 1 layer and 64 units. The last network additionally employs an appearance MLP with 4 layers and 256 units for predicting view-dependent color.

Training Details: Outside of the parameter partitioning scheme proposed in this work, our model architecture matches that of the Zip-NeRF baseline. In experiments on the Zip-NeRF dataset, we partition all parameters in the second proposal network and final density and appearance network and set $p_r = 0.3$, $K_r = 1$. In experiments on the mip-NeRF 360 dataset, only the final network’s parameters are partitioned and we set $p_r = 0.8$, $K_r = 4$.

To maintain a similar memory footprint with the baseline, we reduce the number of multiresolution hash grid entries per partitioned variable by 4×. This ensures that our method uses roughly as much device memory as the baseline, even if the *total* number of parameters is larger.

Comparison to Baseline: On the Zip-NeRF dataset, our method achieves significantly higher quality than the Zip-NeRF baseline at 400,000 steps; see Tab. 1. Qualitatively, this translates to higher geometric and texture detail; see Fig. 3. We further observe a clear, positive relationship between model quality, training time, and number of cells. In contrast, baseline quality does *not* improve with training time. This strongly suggests model capacity, rather than optimization time, to be Zip-NeRF’s limiting factor.

We perform a similar study on a subset of the mip-NeRF 360 dataset, where a similar relationship is lacking; see Tables 2 and 3. While our method modestly outperforms the baseline at 200,000 steps in terms of PSNR and SSIM, a negative relationship between LPIPS and number of cells is evident, as each parameter set receives fewer training iterations as the number of parameter sets grows under a fixed total training budget. This study further suggests that Zip-NeRF’s model capacity is sufficient for medium-sized scenes.

Runtime Analysis: Although our method achieves higher reconstruction quality, the current implementation trains more slowly than its Zip-NeRF counterpart. The reason for this is twofold: first, our method applies each partitioned layer four times to the baseline’s one; second, partitioned parameters are swapped in and out of device memory during training. Based on back-of-the-envelope calculations, we believe the time required for the latter to be largely avoidable with appropriate optimizations. We further believe the number of training steps can be significantly reduced by taking advantage of model parallelism, assigning each cell to a different device. We leave the design and implementation of both of these directions to future work.

5. Conclusion

In this work, we have introduced InterNeRF, a scalable, out-of-core NeRF model architecture for reconstructing large, multi-room scenes. We demonstrated that parameter interpolation is an effective approach for increasing model capacity without a corresponding increase to memory or compute requirements. While our method demonstrates impressive quality, additional work is needed to reduce training time, compare our approach to other submodel-based approaches, test other interpolation schemes, and validate our method on even larger (e.g. city-scale) scenes.

Acknowledgments. We would like to thank Suhani Vora for providing logistical support.

References

- [1] Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M Seitz, and Richard Szeliski. Building rome in a day. *Communications of the ACM*, 54(10):105–112, 2011. 6
- [2] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-NeRF: A multiscale representation for anti-aliasing neural radiance fields. *ICCV*, 2021. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-NeRF 360: Unbounded anti-aliased neural radiance fields. *CVPR*, 2022. 2, 3
- [4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-NeRF: Anti-aliased grid-based neural radiance fields. *ICCV*, 2023. 2, 3, 6
- [5] Daniel Duckworth, Peter Hedman, Christian Reiser, Peter Zhizhin, Jean-François Thibert, Mario Lučić, Richard Szeliski, and Jonathan T. Barron. Smerf: Streamable memory efficient radiance fields for real-time large-scene exploration, 2024. 7
- [6] Mark Kellogg. 3d gaussian splatting for three.js, 2024. <https://github.com/mkkellogg/GaussianSplats3D>. 2, 6
- [7] 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. *ECCV*, 2020. 2, 6
- [8] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *SIGGRAPH*, 2022. 2
- [9] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. DeRF: Decomposed radiance fields. *CVPR*, 2019. 1
- [10] Daniel Rebain, Wei Jiang, Soroosh Yazdani, Ke Li, Kwang Moo Yi, and Andrea Tagliasacchi. Derf: Decomposed radiance fields, 2020. 7
- [11] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. KiloNeRF: Speeding up neural radiance fields with thousands of tiny MLPs. *ICCV*, 2021. 1, 7
- [12] Johannes Lutz Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. *CVPR*, 2016. 3, 6
- [13] Towaki Takikawa, Alex Evans, Jonathan Tremblay, Thomas Müller, Morgan McGuire, Alec Jacobson, and Sanja Fidler. Variable bitrate neural fields. *ACM Transactions on Graphics*, 2022. 7
- [14] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben Mildenhall, Pratul Srinivasan, Jonathan T. Barron, and Henrik Kretzschmar. Block-NeRF: Scalable large scene neural view synthesis. *CVPR*, 2022. 1, 7
- [15] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-NeRF: Scalable construction of large-scale nerfs for virtual fly-throughs. *CVPR*, 2022. 1, 7
- [16] Haithem Turki, Jason Y. Zhang, Francesco Ferroni, and Deva Ramanan. Suds: Scalable urban dynamic scenes, 2023. 7
- [17] Haithem Turki, Michael Zollhöfer, Christian Richardt, and Deva Ramanan. Pynrf: Pyramidal neural radiance fields, 2023. 7
- [18] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and Wenping Wang. F²-nerf: Fast neural radiance field training with free camera trajectories, 2023. 7
- [19] Xiuchao Wu, Jiamin Xu, Xin Zhang, Hujun Bao, Qixing Huang, Yujun Shen, James Tompkin, and Weiwei Xu. ScaNeRF: Scalable bundle-adjusting neural radiance fields for large-scale scene rendering. *ACM Transactions on Graphics (TOG)*, 2023. 1, 7
- [20] Yuanbo Xiangli, Linning Xu, Xingang Pan, Nanxuan Zhao, Anyi Rao, Christian Theobalt, Bo Dai, and Dahua Lin. Bungeenerf: Progressive neural radiance field for extreme multi-scale scene rendering. In *European conference on computer vision*, pages 106–122. Springer, 2022. 6
- [21] Linning Xu, Vasu Agrawal, William Laney, Tony Garcia, Aayush Bansal, Changil Kim, Samuel Rota Bulò, Lorenzo Porzi, Peter Kotschieder, Aljaž Božič, Dahua Lin, Michael Zollhöfer, and Christian Richardt. VR-NeRF: High-fidelity virtualized walkable spaces. *SIGGRAPH Asia*, 2023. 7
- [22] Linning Xu, Yuanbo Xiangli, Sida Peng, Xingang Pan, Nanxuan Zhao, Christian Theobalt, Bo Dai, and Dahua Lin. Grid-guided neural radiance fields for large urban scenes, 2023. 7

InterNeRF: Scaling Radiance Fields via Parameter Interpolation

Supplementary Material

	Network	Levels	Buckets	Features	Interp.	Resol.
Base.	Prop1	6	2^{21}	1	No	512
	Prop2	8	2^{21}	1	No	2048
	Final	10	2^{21}	4	No	8192
Ours	Prop1	6	2^{21}	1	No	512
	Prop2	8	2^{21}	1	No	2048
	Final	10	2^{19}	4	Yes	8192

Table 4. Multiresolution hash grid sizes on the mip-NeRF 360 dataset.

	Network	Levels	Buckets	Features	Interp.	Resol.
Base.	Prop1	6	2^{21}	1	No	512
	Prop2	8	2^{21}	1	No	2048
	Final	10	2^{21}	4	No	8192
Ours	Prop1	7	2^{21}	1	No	1024
	Prop2	9	2^{19}	1	Yes	4096
	Final	11	2^{19}	4	Yes	16384

Table 5. Multiresolution hash grid sizes on the Zip-NeRF dataset.

6. Training Details

Training: We implement our method and the Zip-NeRF baseline by building on the `camp.zipnerf` codebase. Models are trained using sixteen A100s for up to 24 hours. We optimize all models with a batch size of 2^{16} rays using the Adam optimizer. Unless otherwise stated, hyperparameters match those described in the Zip-NeRF paper [4]. We employ an exponentially decaying learning rate schedule with an initial learning rate of 1e-2 and a final learning rate of 1e-3 or 1e-4 on the Zip-NeRF and mip-NeRF 360 datasets, respectively. We further introduce a linear learning rate warm-up schedule of 2,500 steps starting with an initial learning rate of 1e-8. We apply L_2 regularization to all hash grid parameters with a weight of 0.001 for proposal networks and 0.1 for density and appearance networks on the Zip-NeRF dataset and 0.1 for all networks on the mip-NeRF 360 dataset. We further force volumetric rendering weights along each camera ray to sum to unity.

Model Architecture: All models, including the Zip-NeRF baseline and InterNeRF, are composed of three networks: two proposal networks and one appearance and density network. Each network, in turn, is composed of up to three components: a multiresolution hash grid, a Geometry MLP, and potentially an Appearance MLP. All three networks contain the first two components; only the last contains an Appearance MLP. The shape of these MLPs is identical for all methods as described in the main text.

The majority of each method’s parameters lie in their multiresolution hash grid. We describe the number of hash resolution levels, number of entries per level, number of features per entry, and spatial resolution of these multiresolution grids in Tables 4 and 5. We further indicate which multiresolution grids are spatially-partitioned and interpolated for our method.

Cell Swapping: We train our method by swapping INTERPOLATED model parameters in and out of device memory as needed. As described in Sec. 3, scene grid cells are op-

timized one at a time. We optimize each cell k for $2 \times N_k$ iterations, where N_k is the number of training cameras allocated to this cell. When a new grid cell is chosen, SHARED parameters are left in memory while the appropriate INTERPOLATED are swapped in. To order cells, we assign each cell a linearized integer identifier based on its position in the 2D scene grid and loop over cells according to this identifier. When loading a new cell for training, INTERPOLATED variables along with associated Adam statistics are loaded into memory.

Dataset: Unlike the majority of prior work, we use high-resolution versions of the mip-NeRF 360 and Zip-NeRF datasets with the original lens distortion. This enables a deeper exploration of model capacity – a key goal of this work – but prevents metrics from being directly comparable to prior publications. The use of distorted photos further limits our ability to compare to rasterization-based methods such as 3D Gaussian Splatting [6]. Concretely, we use photos at double the resolution used in prior work. For the mip-NeRF 360 dataset, this means full-resolution photos for indoor scenes and $2 \times$ downsampled photos for outdoor scenes. For the Zip-NeRF dataset, we use $2 \times$ downsampled photos for all scenes except for BERLIN, where we use $4 \times$ downsampling.

7. Related Work

Large-Scale 3D Reconstruction: There is a large body of work tackling large scene reconstruction. This includes a number of classic works that apply structure from motion to extremely large photo collections [1, 12]. Ever since neural radiance fields [7] emerged as the dominant paradigm for novel view synthesis, a number of works have sought to also scale NeRF to large scenes. BungeeNeRF [20] combines satellite and ground level images in a single NeRF by appending residual blocks to represent progressively finer scales. Each block has its own output head that predicts residual color and density, but earlier blocks are only trained on coarser scale (aerial) images. Similar ideas for rendering

at high resolution in a progressive manner are explored in variable bitrate neural fields and PyNeRF [13, 17].

Grid-Guided NeRFs [22] initially trains instant-NGP with a small MLP and then later incorporates a large view-dependent MLP, an approach that helps the model perform better on large urban scenes. F2-NeRF [18] proposes adaptive space warping by finding local perspective warps that shrink regions far from any camera (a generalization of the NDC transform to multiple cameras), which helps to allocate model capacity more efficiently in scenes with long camera trajectories. VR-NeRF [21] uses a customized camera rig to capture high resolution HDR footage that can be used to train a large-scale NeRF.

Submodel NeRFs: Several works train independent NeRF submodels that are each responsible for representing a subset of the scene, and aggregate their predictions at inference time. Some works were motivated by efficiency or composition [10, 11]. In Block-NeRF [14], street view data is represented by submodels centered at street intersections, each trained on cameras within a fixed radius. Target views are rendered by compositing whichever submodels have visibility of the region, as estimated by an auxiliary visibility network. Mega-NeRF [15] partitions a scene using an octree, and trains each cell only on rays that pass through it. Scalable Urban Dynamic Scenes [16] builds on this work by using LIDAR data to prune rays that do not terminate within the cell. ScaNeRF [19] uses a similar tile-based submodel partition, and additionally performs bundle adjustment within a distributed parallel training framework. Streamable MERF [5] distills a large NeRF into several submodels that use a memory-efficient triplane representation. The work also features a deferred MLP whose parameters are interpolated based on the distance of the camera origin from the submodel centers, similar to the way parameters are interpolated in our work.

Overall, we find a strong positive correlation between model capacity and texture detail on multi-room scenes as demonstrated in Fig. 5. Our method consistently provides crisp detail on textured surfaces including painting, carpets, and curtains. The same cannot reliably be said, however, for the mip-NeRF 360 dataset. In the majority of scenes, our method is indistinguishable from the Zip-NeRF baseline as in Fig. 4.

8. Additional Results

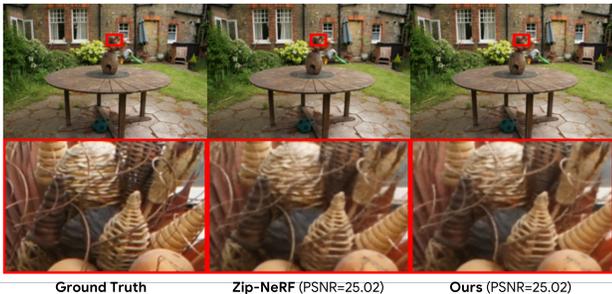


Figure 4. Renders of the GARDEN scene in the mip-NeRF 360 dataset using Zip-NeRF and InterNeRF with a 3×3 parameter grid.

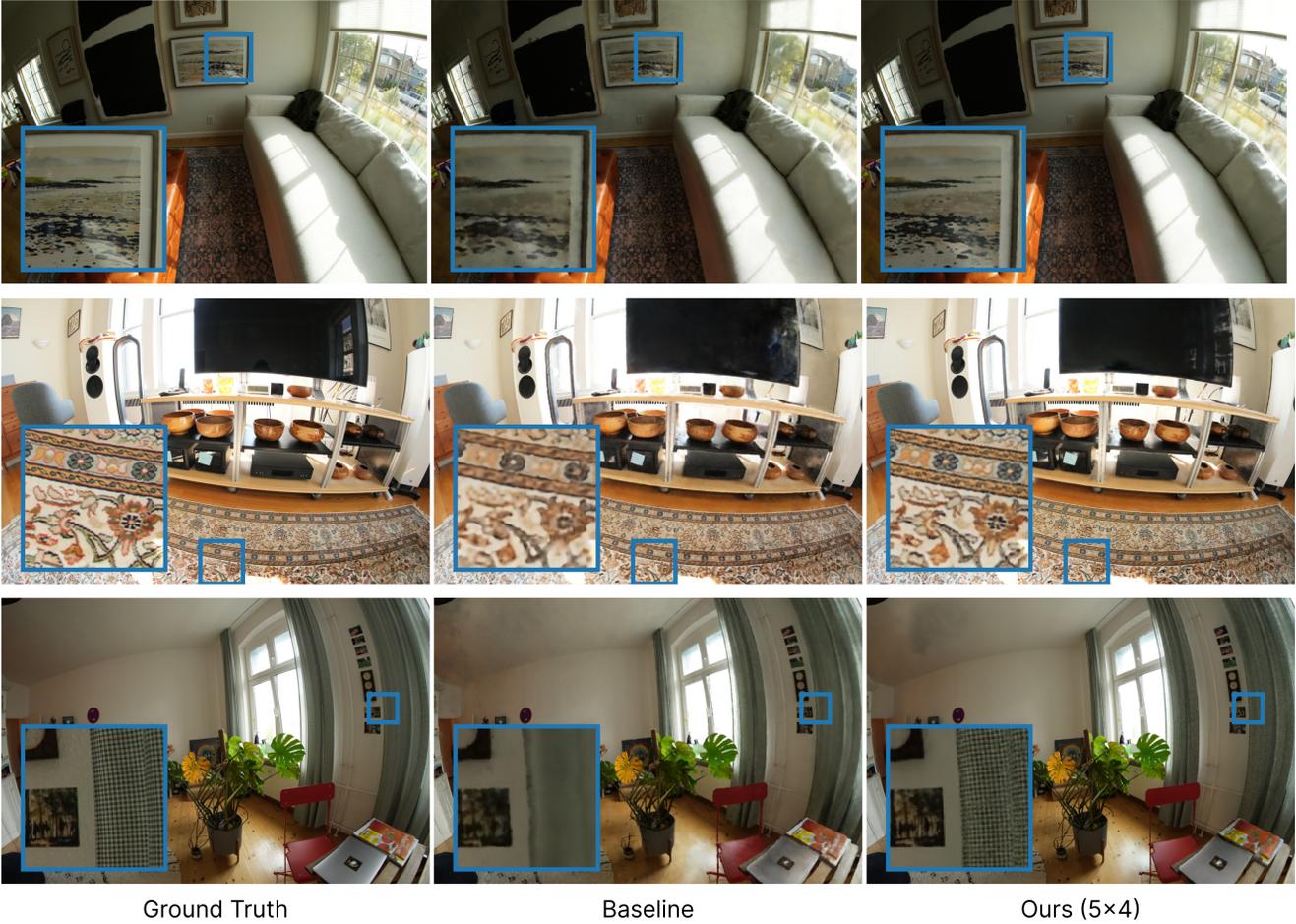


Figure 5. Additional qualitative results on the Zip-NeRF dataset. Our method consistently reconstructs centimeter-level texture detail in large, multi-room scenes.