

SCALAR-NeRF: SCALable LARge-scale Neural Radiance Fields for Scene Reconstruction

Yu Chen

Gim Hee Lee

Department of Computer Science, National University of Singapore

chenyu@comp.nus.edu.sg

gimhee.lee@nus.edu.sg

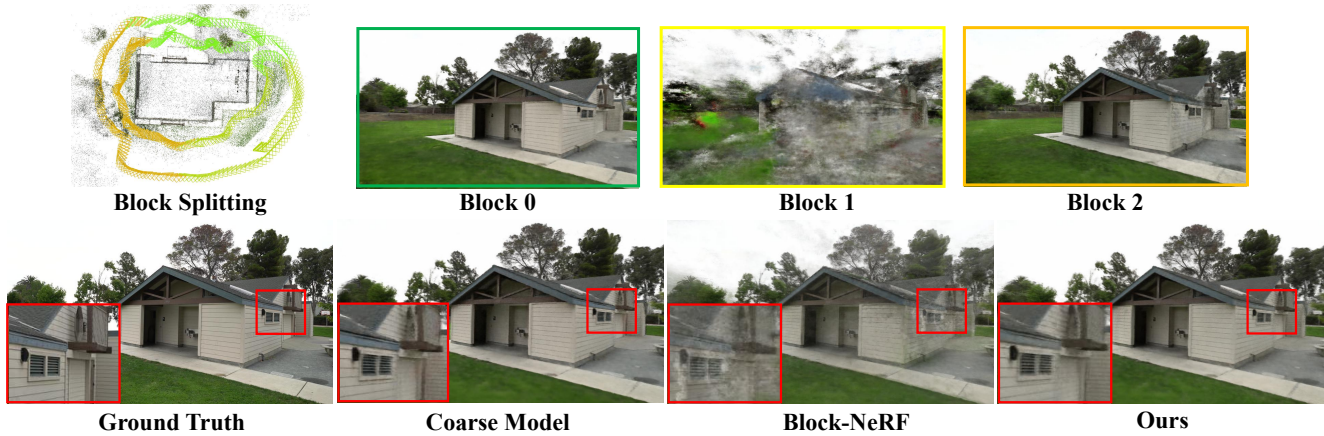


Figure 1. **Comparison of our method to Block-NeRF.** Our method adopts a coarse-to-fine framework to train NeRF models under large-scale scenes, which effectively avoids the ‘foggy’ issue in Block-NeRF.

Abstract

In this work, we introduce *SCALAR-NeRF*, a novel framework tailored for scalable large-scale neural scene reconstruction. We structure the neural representation as an encoder-decoder architecture, where the encoder processes 3D point coordinates to produce encoded features, and the decoder generates geometric values that include volume densities of signed distances and colors. Our approach first trains a coarse global model on the entire image dataset. Subsequently, we partition the images into smaller blocks using *KMeans* with each block being modeled by a dedicated local model. We enhance the overlapping regions across different blocks by scaling up the bounding boxes of each local block. Notably, the decoder from the global model is shared across distinct blocks and therefore promoting alignment in the feature space of local encoders. We propose an effective and efficient methodology to fuse the outputs from these local models to attain the final reconstruction. Employing this refined coarse-to-fine strategy, our method outperforms state-of-the-art NeRF methods and demonstrates scalability for large-scale scene reconstruction. The code will be available on our [project page](#).

1. Introduction

Over the past three decades, 3D scene reconstruction remains an active research field with useful applications spanning augmented reality, real-world simulation, localization, etc. The advent of neural scene reconstruction techniques [32, 49] has revolutionized the digitalization of 3D scenes, unveiling a realm beyond the limitations of traditional photogrammetry tools [40]. These neural methods exhibit superior resilience to non-Lambertian effects, appearance changes [29], and dynamic scenes [36] compared to their traditional counterparts.

Recent advancements have significantly improved the efficiency of training neural scene representations. While earlier approaches require long training time on a vanilla NeRF [32], innovative techniques such as Instant-NGP [33] and TensoRF [7] have remarkably condensed this training time to mere minutes. Nonetheless, the majority of these developments primarily focus on small-scale object-centric scenes, leaving the frontier of neural reconstruction on large-scale scenes largely unexplored.

The pursuit of neural reconstruction in large-scale scenes is not an easy endeavour, where significant challenges arise from several factors:

- **Parameter Overload.** The sheer volume of parameters

required for updating poses a significant obstacle. Although augmenting model parameters might enhance reconstruction quality, it becomes challenging to train and manage such large models on a single GPU.

- **Extended Training Duration.** Larger model capacities demand extended training times to capture high-fidelity scene details, thereby posing a time-intensive challenge.
- **Incomplete Scene Details.** Mere augmentation of model parameters may still result in missing scene details, highlighting the limitations of this simplistic approach.

To overcome inherent limitations and extend neural scene reconstruction to city-scale scenes, recent works such as Block-NeRF [44] and Mega-NeRF [45] have embraced a divide-and-conquer approach. While Block-NeRF adeptly addresses the challenge of maintaining appearance consistency across different blocks, noteworthy challenges persist, particularly when block centroids are in close proximity. In such cases, Block-NeRF tends to generate images with a foggy appearance, attributed to its straightforward inverse distance weighting function for image blending. On the other hand, Mega-NeRF primarily focuses on utilizing aerial images for applications such as rescue searches, presenting its unique set of constraints and limitations.

In this paper, we introduce a novel coarse-to-fine framework for reconstructing large-scale scenes using NeRF. Our method leverages a multi-resolution hashing table as the feature encoder. Although the multi-resolution hash table is widely adopted in neural scene reconstruction, we observe that simply increasing the hash table size not only slows down training but also renders it infeasible to fit into a single GPU during training.

The core insight of our method is that a large hash table cannot be accommodated by a non-commercial GPU, while a small hash table has a higher probability of hash collisions. Splitting scenes into smaller blocks enables training local models with finer details without the necessity of increasing the hash table size. Initially, we train a rapid, coarse global model with all training images. Despite its expedited training, this coarse model retains considerable reconstruction quality. Subsequently, we partition the training images into smaller blocks, each reconstructed using an individual local model. We further enhance the intersection areas of local blocks by expanding the bounding boxes that cover the images in the current block. To improve the regularization of the hash encoder in each local model, we share the decoder from the coarse model across all local models. The shared decoder can either be fixed for faster training speed or fine-tuned for better reconstruction quality. Unlike Block-NeRF, which blends images using the inverse distance weighting function, we synthesize images from each block and fuse them with a reference image rendered from the coarse model. Leveraging the *coarse-global to fine-local* strategy, our method with a hash table size of

2^{19} surpasses Block-NeRF and Instant-NGP with a hash table size of 2^{21} .

We train and evaluate our model on outdoor large-scale scenes. Experiments demonstrate its substantial performance, surpassing existing state-of-the-art methods such as Block-NeRF. The improvements are evident in the scene reconstruction showcased in Fig. 1. The main contributions of our work are:

- We propose a novel framework for reconstructing large-scale scenes.
- The proposed method enables training NeRF models under large-scale scenes on a single GPU without compromising reconstruction quality.
- Extensive experiments on the large-scale Tanks and Temples dataset show that our method surpasses existing NeRF methods.

2. Related Work

Large-Scale 3D Reconstruction. Snavely *et al.* [41] pioneered a system enabling virtual visits to global attractions. They utilized Structure-from-Motion (SfM) [40] and key-point extraction with SIFT [27] to reconstruct sparse scene structures. Image similarity searching methods like vocabulary trees [34] and NetVLAD [1] were applied to enhance reconstruction efficiency by removing unnecessary matching pairs. The works by Zhu *et al.* [58, 59] and Chen *et al.* [10, 11] adopted a divide-and-conquer strategy for scalability under large-scale scenes. The improvement of reconstruction quality and efficiency in large-scale scenes has been significantly influenced by Multi-View Stereo (MVS) techniques. To recover the scene details as much as possible, MVS [14, 15] is used to densify scene structures by plane sweep [16] or patch-match [2]. Recently, Yao *et al.* [54] leverages 3D convolutions to learn the depth probability map from the cost volume constructed by plane sweep. Later works [28, 47, 48, 55] also follow a similar paradigm to improve the reconstruction quality and efficiency of MVS.

Novel View Synthesis. The field of novel view synthesis has seen significant advancements, notably with neural radiance fields [32] enabling rendering from novel viewpoints with encoded frequency features [43]. To improve the rendering efficiency, NSVF [25] represents scenes into sparse voxels, where the voxels are gradually updated by checking the volume densities. PlenOctrees [56] also uses sparse voxels for scene representation. To improve the training speed and rendering efficiency, the multi-layer perceptron (MLP) is converted to volume densities and spherical harmonics that are stored on the octree leaves. Plenoxels [13] found that the MLP is not necessary, and instead uses trilinear interpolation to improve the continuities of the scalar values stored on the voxel corners. Instant-NGP [33] encoded the features into a multi-resolution hash table, where

the hash collision is implicitly handled during optimization. TensorRF [7] uses CP-decomposition or VM-decomposition to encode scenes into three orthogonal axes and planes. The low-rank decomposition results in highly compact representations. Gaussian Splatting [20] initializes 3D Gaussians from a set of sparse point clouds, the 3D Gaussians are used as explicit scene representation and dynamically merged and split during training. The splatting operation can be performed very fast through rasterization, which enables real-time novel view synthesis. Other methods also focus on the generalizability of NeRF [6, 19, 26, 42, 50], bundle-adjusting camera poses and NeRF [9, 24, 30], removing floaters [18, 39], and leveraging sparse or dense depth to supervise the training of NeRF [12, 38, 52], *etc.*

Large-Scale NeRF. To enhance the representation ability of NeRF on outdoor scenes, Mip-NeRF [3] proposed to use Gaussian to approximate the cone sampling, the integrated positional encodings are therefore scale-aware and can be used to address the aliasing issue of NeRF. Mip-NeRF360 [4] further uses space contraction to model unbounded scenes. Zip-NeRF [5] adopted a hexagonal sampling strategy to handle the aliasing issue for Instant-NGP [33]. NeRF-W [29] models the scene changes and transient effects by attaching an appearance encoding for each image. Block-NeRF [44] and Mega-NeRF [45] adopt a similar divide-and-conquer strategy as [10, 59], where smaller blocks can be reconstructed in parallel. To this end, Block-NeRF focus on fixing the appearance inconsistency issue between different blocks, while Mega-NeRF aims at encouraging the sparsity of the network under aerial scenes. Urban-NeRF [37] leverages lidar points to supervise the depth of NeRF in outdoor scenes. SUDS [46] further extended Mega-NeRF into dynamic scenes. Different to previous large-scale NeRF methods, Switch-NeRF [31] uses a switch transformer that learns to assign rays to different blocks during training. Grid-NeRF [53] designed a two-branch network architecture, where the NeRF branch can encourage the feature plane [7] branch recover more scene details under large-scale scenes. However, the two-branch training scheme is trivial and needs a long time to train. NeRF2NeRF [17] and DReg-NeRF [8] assumes images are only available during training in each block, and they propose methods to register NeRF blocks together.

3. Our Method

Fig. 2 shows an illustration of our network architecture. Given the full set of training images, we first train a coarse global model. Subsequently, we split the images into different blocks. We further create overlapping regions across different blocks by scaling up the bounding boxes which cover all images in local blocks. We then share the decoder from the global model across all local blocks, and finetun-

ing it to better condition the hash encodings in local blocks. The final rendered images are fused from different blocks with the global model as guidance.

3.1. Preliminaries

Neural Radiance Fields (NeRF). Neural Radiance Fields (NeRF), proposed by Mildenhall *et al.* [32] in 2020, revolutionized 3D scene reconstruction and novel view synthesis. NeRF represents scenes as a continuous function, estimating the volumetric scene representation by predicting the volume density and view-dependent emitted radiance at any given point. As introduced by NeRF [32], the fundamental equation for volume rendering can be defined as:

$$C(x) = \int_{t_{\text{near}}}^{t_{\text{far}}} T(t) \cdot \sigma(t) \cdot L_i(t) dt, \quad (1)$$

where $C(x)$ denotes the color of a pixel, $\sigma(t)$ is the volume density, $L_i(t)$ is the emitted radiance, and t_{near} and t_{far} indicate the start and end points of the ray traversing the scene, respectively. $T(t)$ is the accumulated transmittance:

$$T(t) = \exp\left(-\int_{t_{\text{near}}}^t \sigma(s) ds\right). \quad (2)$$

Instant Neural Graphics Primitives for NeRF. As an extension of NeRF introduced by Müller *et al.* [33], Instant Neural Graphics Primitives (InstantNGP) offers a solution to the computational challenges associated with NeRF where training time is significantly reduced from days to minutes by encoding scene representations. By addressing hash collisions implicitly, InstantNGP enables faster and more memory-efficient neural scene reconstruction, offering a practical solution for large-scale scene reconstruction with reduced computational demands.

In this work, we utilize InstantNGP as our NeRF backbone, where the multi-resolution hash table is formulated as an encoder:

$$\mathbf{F} = \text{Encoder}(\mathbf{x}), \quad (3)$$

where \mathbf{F} is the feature embedding concatenated from the multi-resolution hash table $\mathbf{F} = [\mathbf{f}_1, \dots, \mathbf{f}_L] \in \mathbb{R}^{L \cdot F}$, L is the number of levels, and F is the feature dimension. The embedding is then decoded into volume density σ and per-point color \mathbf{c} by a shallow layer MLP:

$$\sigma, \mathbf{c} = \text{Decoder}(\mathbf{F}). \quad (4)$$

For unbounded scenes, we contract the 3D point \mathbf{x} using the space contraction method of Mip-NeRF360 [4]:

$$\text{contract}(\mathbf{x}) = \begin{cases} \mathbf{x}, & \|\mathbf{x}\| \leq 1, \\ \left(2 - \frac{1}{\|\mathbf{x}\|}\right) \begin{pmatrix} \mathbf{x} \\ \mathbf{x} \end{pmatrix}, & \|\mathbf{x}\| > 1. \end{cases} \quad (5)$$

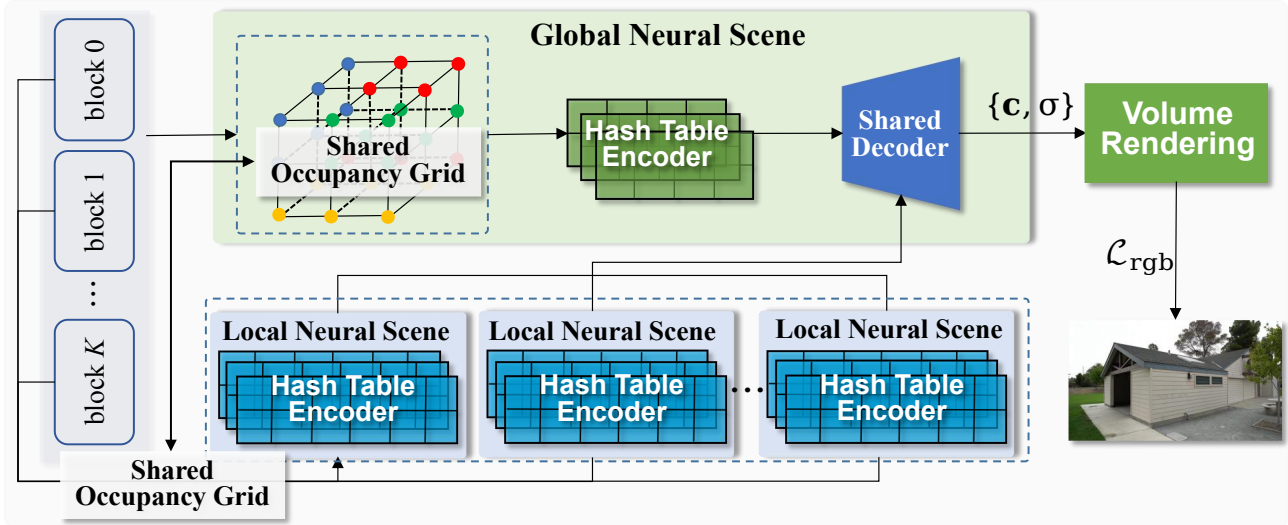


Figure 2. **Network architecture of our method.** Given the full set of training images, we first train a coarse global model. Subsequently, we split images into different blocks. We further create overlapping regions across different blocks by scaling up the bounding boxes which cover all images in local blocks. We then share the decoder from the global model across all local blocks, and finetuning it to better condition the hash encodings in local blocks. The final rendered images are fused by a ‘winner-takes-all’ strategy.

3.2. Coarse Global Reconstruction

Our initial step involves training a coarse global model to represent the entire scene. The hash function maps a point to its hash entry is defined as:

$$h(\mathbf{x}) = \left(\bigoplus_{i=1}^d x_i \pi_i \right) \bmod T, \quad (6)$$

where T is total number of hash entries at each level, \oplus denotes the bit-wise XOR operation, and π_i are large unique prime numbers. The global model which we termed as the “coarse model” shares the same hash table size as the local models. However, since a whole scene contains more points \mathbf{x} than a part of the scene and the hash table size T remains constant, hash collisions are more likely to occur in the global model compared to a local model with the same hash table size. While Instant-NGP implicitly handles hash collisions during training, collided hash entries in the global model may take a weighted mean of gradients from different parts of the scene. This equality in contribution from various parts can degrade the quality of reconstruction in the final model.

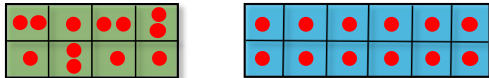


Figure 3. **An illustration of hash collision.** Green table: a small hash table with 8 entries; Blue table: a larger hash table with 12 entries. Red circle: 3D points that are hashed into the bins.

We observe that enlarging the size of the hash table is impractical for two main reasons:

- **Larger model requires more memory and longer training time.** For implicit representations like MLP, a

larger model demands more memory during training due to updates to all parameters during backpropagation. This increase in computational demands slows down training. For InstantNGP, although only the feature vectors of local voxels are updated, a larger memory is still required. Fig. 3 illustrates why training a larger hash table is slower when given a fixed number of points. Suppose we have twelve 3D points $\{\mathbf{x}_i \mid i \in [1, N]\}$, a small hash table \mathcal{H}_s , and a larger hash table \mathcal{H}_l . For \mathcal{H}_l , gradients must be computed for all bins $[\frac{\partial \mathcal{L}}{\partial \mathbf{F}_1}, \dots, \frac{\partial \mathcal{L}}{\partial \mathbf{F}_{12}}]$. However, for \mathcal{H}_s , only eight bins need computing gradients during backpropagation $[\frac{\partial \mathcal{L}}{\partial \mathbf{F}_1}, \dots, \frac{\partial \mathcal{L}}{\partial \mathbf{F}_8}]$, as four bins contain more than one point due to hash collisions.

- **Larger model does not always translate to better reconstruction quality.** This is evident in Fig. 5, where undeterministic hash collisions can result in two parts contributing equally to the same voxels.

3.3. Fine Local Reconstruction

Block Splitting. Given the complete set of training images $\{\mathbf{I}_1, \dots, \mathbf{I}_N\}$ and their corresponding camera poses $\{\mathbf{P}_1, \dots, \mathbf{P}_N\}$. traditional distributed Structure from Motion (SfM) methods [10, 59] often adopt a recursive strategy to construct intersected blocks. In detail, a view graph is constructed with images as nodes and the number of correspondences as edge weight. The graph is then split into disconnected components, and images in the current component are inserted into other components based on the number of correspondences. However, we find this strategy unsuitable for our purposes. This is because novel view synthesis in NeRF relies on photometric information, but SfM primarily relies on feature-metric correspondences. Conse-

quently, the block-splitting method in SfM might yield sub-optimal results, especially in texture-less regions. Additionally, distributed SfM necessitates splitting images before computing camera poses, while NeRF already has known camera poses and thus enabling us to leverage these priors for our block-splitting method.

We employ KMeans to cluster images close to each other into the same block: $L_{I_i} = \text{KMeans}(\mathbf{t}_i)$, where \mathbf{t}_i represents camera centers, L_{I_i} is the block ID of image i . While Mega-NeRF [45] also uses KMeans to split blocks, the performance is often unsatisfactory due to the lack of overlapping regions between blocks. To this end, we compute axis-aligned bounding boxes AABB_k using the camera poses in each block. We then scale up each AABB using a scale factor vector $\mathbf{s}_{\text{AABB}} = [s_x, s_y, s_z]^\top$, where s_x, s_y, s_z are the scale factor along the x, y, z axis. By default, $s_x = s_y = s_z$. Next, we regroup images into different blocks, where images falling into the same AABB belong to the same block. This approach introduces overlapping regions between blocks since the AABBs intersect. By controlling \mathbf{s}_{AABB} , we can manage the overlapping ratio across blocks: 1) Blocks have no overlapping regions when $s = 1.0$; 2) Blocks encompass all training images when $s \rightarrow \frac{\text{AABB}_{\text{whole}}}{\text{AABB}_k}$, where $\text{AABB}_{\text{whole}}$ denotes the AABB for the whole training images. Our block-splitting algorithm is summarized in Alg. 1.

Algorithm 1 Block Splitting Algorithm

Require: Camera poses $\mathcal{P} = \{\mathbf{P}_1, \dots, \mathbf{P}_N\}$; Bounding-box scale factor vector $\mathbf{s}_{\text{AABB}} = [s_x, s_y, s_z]^\top$; number of blocks K .

Ensure: Intersected blocks $\mathcal{B} = \{\mathbf{B}_k = \{\mathbf{I}_{i,k}\} \mid k \in [1, K]\}$.

- 1: Obtain no overlapping blocks by $\mathcal{B} \rightarrow \text{KMeans}(\mathcal{P}, K)$.
- 2: **for** $k \in [1, K]$ **do**
- 3: Compute the bounding box for block k :

$$\text{AABB}_k = [\mathbf{A}, \mathbf{B}] = [\min(\mathcal{P}_k), \max(\mathcal{P}_k)]. \quad (7)$$

- 4: Compute the midpoint of the diagonal for AABB_k : $\mathbf{C} = (\mathbf{A} + \mathbf{B})/2$.
- 5: Compute the ray direction of the diagonal $\mathbf{r} = (\mathbf{A} - \mathbf{B}) / (\|\mathbf{A} - \mathbf{B}\|)$, half of the length of the diagonal $\text{len} = \|\mathbf{B} - \mathbf{A}\|/2$.
- 6: Recompute \mathbf{A} and \mathbf{B} by:

$$\begin{aligned} \mathbf{A} &= \mathbf{C} + \mathbf{s}_{\text{AABB}} \cdot \mathbf{r} \cdot \text{len}, \\ \mathbf{B} &= \mathbf{C} - \mathbf{s}_{\text{AABB}} \cdot \mathbf{r} \cdot \text{len}. \end{aligned} \quad (8)$$

- 7: Update block $\mathbf{B}_k \rightarrow \mathbf{B}_k + \mathbf{I}_i, \forall \mathbf{B} \leq \mathbf{P}_i \leq \mathbf{A}$.
-

Global Conditioned Local NeRF. With the intersected blocks now constructed, we proceed to train a local NeRF model for each block individually. Benefiting from fewer

hash collisions, the local model can capture a higher fidelity of the scene block, even with a hash encoder featuring the same hash table size as the coarse global model. To further enhance the reconstruction quality, we adopt a strategy where the decoder from the coarse global model is shared across all local models.

The shared decoder is trained using all available training images, allowing it to effectively condition the local hash encoder and align the feature space of hash encodings in different local models. To optimize training speed and reduce memory footprint, each local model receives a copy of the shared decoder, and fine-tunes its own copy during training. This approach leverages the comprehensive knowledge of the shared encoder while allowing each local model to adapt to the specific details within its corresponding block.

3.4. Final Image Blending

Blending images from different blocks poses a non-trivial challenge as there are no explicit priors available to guide the weighting of image colors. Traditional methods such as the inverse distance weighting (IDW) scheme employed by Block-NeRF often lead to undesired outcomes of foggy or blurry images.

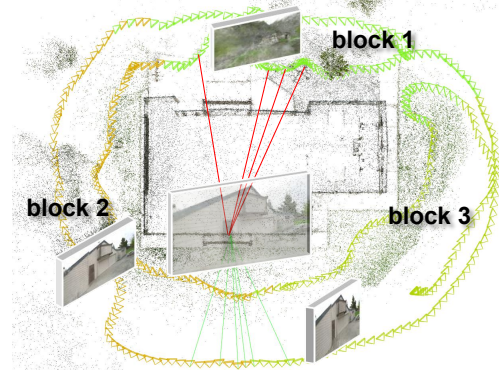


Figure 4. **Using IDW tends to generate blurry images when blocks are close to each other.** Left: The scene is split into three blocks, while the point on the right-top part of the barn is only visible by two blocks. Right: Block 1 does not observe that point. However, the three blocks are close to each other. IDW therefore gives approximately the same weight to all blocks when blending images, which includes outliers from Block 1.

Inverse Distance Weighting. The IDW function to blend images from different blocks is defined as follows:

$$\mathbf{I} = \sum_k^K \omega_k \mathbf{I}_k, \quad \omega_k = d_k^{-\gamma}, \quad (9)$$

where d_k represents the distance between the query point and the block center, and γ is a hyperparameter (default to $\gamma = 1$ in our experiment). Despite its widespread use, the IDW approach tends to generate blurry images when blocks are in close proximity, as illustrated in Figure 4.

Image Blending with Global Guidance. An alternative solution involves computing uncertainties for the rendered color of each pixel and discarding results from blocks with significant uncertainties. However, existing uncertainty computation methods designed for NeRF, *e.g.*, NeRF-W [29] and Bayes’ Rays [18], are found to be insufficiently reliable. Fortunately, a simple yet effective approach involves leveraging guidance from the coarse global model to fuse high-quality images. This is achieved by comparing the color difference between images rendered from local models and the coarse model. The coarse model is employed to compute the color for a point as well as local blocks. The pixel with the minimum color difference is then selected as:

$$\mathbf{c} = \min(\|\mathbf{c}_k - \mathbf{c}_{\text{global}}\|). \quad (10)$$

It is noteworthy that to maintain efficient inference, not all local models need to be considered in the comparison with the global model. As suggested by Block-NeRF, filtering out blocks that are highly unlikely to be observed in the current region can be achieved by training a visibility field. However, we did not implement this filtering mechanism in our experiments since our blocks are sufficiently close. We empirically affirm that this choice does not compromise the contribution and practicality of our method.

4. Experiments

Implementations. We use PyTorch [35] for our implementation. Training is performed using the Adam optimizer [21] for 200,000 iterations. We compare our approach with baseline methods, including vanilla NeRF [32], Instant-NGP [33], and our reimplementation of Block-NeRF [44] since the original code for Block-NeRF is not publicly available. The learning rate is set to $5e-4$ for vanilla NeRF and $1e-3$ for all other methods. Notably, the learning rate for fine-tuning the shared decoder in our method is set to $5e-5$. A warm-up phase of 5,000 iterations is employed with the learning rate linearly increasing during this period, followed by exponential decay with a rate of 0.1. To enhance training efficiency, we incorporate the occupancy grid sampler from NeRFacc [23] to skip empty space in our method. The occupancy grid sampler is also applied to all other methods for a fair comparison. It is important to note that the fine branch of the original vanilla NeRF is not enabled in our experiments due to memory constraints and increased training time. All experiments are conducted on an NVIDIA A5000 GPU with 24GB memory. By default, the hash table size for both our coarse global model and local models is set to 2^{19} . The hash tables have 16 levels of resolution, and each level features dimensions of 4. Our decoder has a depth of 5 with a width of 128. These hyperparameters are selected based on comprehensive experimentation and provide a good trade-off between

performance and resource utilization.

Datasets. Our method is trained and evaluated on different large-scale unbounded scenes of the tanks-and-temples dataset [22]. This dataset poses significant challenges as it encompasses diverse scenes with varying lighting conditions, and appearance changes (such as overexposure and shadows on the ground), and contains 300 to 1,200 images per scene. To split the dataset for training and validation, we adopt a strategy of holding out every 20 images for validation and utilizing the remaining images for training. To model appearance changes, we follow the approach of NeRF-W [29] by associating each image with an appearance embedding. The appearance embedding dimension is set to 8 for all methods. During validation, the appearance embeddings for the validation images are unknown. To address this, we calculate the mean of the appearance embeddings for all training images and employ it as the appearance embedding for the validation images. This approach ensures a consistent representation of appearance across the dataset during validation.

Results. We employ PSNR, SSIM [51] and LPIPS [57] as metrics for novel view synthesis. Quantitative results are presented in Tab. 1 with qualitative comparisons shown in Fig. 6. From Tab. 1, it is evident that our method excels in terms of SSIM and achieves the second-best result in terms of PSNR. Although vanilla NeRF [32] attains the best PSNR, it consistently produces foggy images (*cf.* the second column in Fig. 6). For InstantNGP, we provide its results with different hash table sizes of 19 and 21. From Tab. 1, we can see that NGP_{21} achieves the best results in LPIPS and the second-best in SSIM. The reconstruction quality decreases with a reduced hash table size (*cf.* NGP_{19} in Tab. 1). Importantly, a larger hash table size does not consistently yield superior results. For the ‘Courthouse’ scene, NGP_{21} obtained the worst result in PSNR and performs worse than NGP_{19} in terms of SSIM and LPIPS. In Fig. 5, we visualize a region where vanilla NeRF and NGP_{21} inaccurately reconstructed the corner of the building. As expected, Block-NeRF consistently produces hazy images for almost all scenes in the tanks-and-temples dataset. Conversely, our method can avoid this issue by leveraging the coarse global model to discard poorly rendered results in the fusion stage. Our method can render higher-quality images than NGP_{21} in details even though the hash table size is smaller than NGP_{21} in detailed areas, even with a smaller hash table size, as observed in the zoomed-in text regions of scenes ‘Caterpillar’, ‘Family’, and ‘Truck’ in Fig. 6. Overall, our method demonstrates superior scalability in both quantitative and qualitative evaluations compared to other methods.

Ablation Study of Network Architecture. We present the ablation study of our network architecture in Tab. 2. The results are averaged from the 8 scenes of the tanks-and-



Figure 5. **Failure cases of Instant-NGP [33] on tanks-and-temple dataset [22].** From left to right are ground-truth image, vanilla NeRF [32], Instant-NGP [33] with hash table size 2^{19} , InstantNGP [33] with hash table size 2^{21} , Block-NeRF [44], our Scalar-NeRF.

Scenes	PSNR \uparrow					SSIM \uparrow					LPIPS \downarrow				
	NeRF	NGP ₁₉	NGP ₂₁	Block-NeRF	Ours	NeRF	NGP ₁₉	NGP ₂₁	Block-NeRF	Ours	NeRF	NGP ₁₉	NGP ₂₁	Block-NeRF	Ours
Barn	22.415	20.422	22.029	19.934	21.921	0.616	0.570	0.638	0.631	0.639	0.470	0.396	0.334	0.361	0.378
Caterpillar	19.955	17.220	18.027	18.117	18.280	0.503	0.506	0.516	0.516	0.541	0.540	0.413	0.381	0.399	0.420
Courthouse	13.073	15.088	12.586	15.010	15.990	0.374	0.504	0.428	0.486	0.500	0.708	0.553	0.660	0.590	0.534
Family	18.832	17.606	17.197	16.618	17.838	0.542	0.516	0.506	0.485	0.471	0.426	0.331	0.307	0.362	0.380
Horse	20.470	17.632	17.917	17.656	18.424	0.686	0.633	0.622	0.635	0.648	0.333	0.294	0.273	0.329	0.338
Ignatius	18.217	15.709	17.285	16.613	17.498	0.431	0.375	0.481	0.367	0.458	0.548	0.416	0.357	0.380	0.405
Lighthouse	17.386	15.369	16.476	16.080	16.280	0.571	0.536	0.558	0.569	0.548	0.517	0.506	0.479	0.494	0.537
Truck	17.788	18.562	18.288	17.722	17.984	0.528	0.575	0.569	0.556	0.527	0.470	0.356	0.310	0.341	0.387
Average	18.517	17.201	17.476	17.219	18.027	0.531	0.527	0.540	0.531	0.542	0.512	0.408	0.388	0.407	0.422

Table 1. **Quantitative results of novel view synthesis on tanks-and-temple [22] dataset.** NGP₁₉ and NGP₂₁ denote InstantNGP [33] with hash table size 2^{19} and 2^{21} , respectively. \bullet \bullet and \bullet respectively denotes the **first**, **second**, and **third**-best results.

temples dataset. The AABB scale factor is $[1.0, 1.0, 1.0]$. Specifically, ‘w.o. SD’ denotes our method without the shared decoder, ‘w.o. FT’ denotes our method without finetuning the shared decoder for each local model, ‘Full (idw)’ denotes our method with inverse distance weighting function for final fusion, ‘Full (global)’ denotes our method with the coarse global model for final fusion. As observed, the shared decoder plays a critical role in our network architecture. Moreover, the local model without finetuning the shared decoder can already generate satisfactory results. It suggests that our network architecture has the potential to train larger models by increasing the hash table size and fixing the shared decoder in each local model, eliminating the need to save the computational graph for the parameters of the shared decoder during training. The results indicate that our method with the IDW function performs worse than using the results from the coarse global model, revealing that the IDW function is not suitable when blocks are close to each other. Finally, our method with the full pipeline achieves the best results.

	w.o. SD	w.o. FT	Full (idw)	Full (global)
PSNR \uparrow	15.384	17.069	17.167	17.512 \bullet
SSIM \uparrow	0.490	0.537	0.549	0.565 \bullet
LPIPS \downarrow	0.521	0.451	0.437	0.435 \bullet

Table 2. **Ablation studies of our network architecture.** ‘w.o. SD’ denotes our method without shared decoder, ‘w.o. FT’ denotes our method without finetuning the shared decoder for each local model, ‘Full (idw)’ denotes our method with inverse distance weighting function for final fusion. ‘Full (global)’ denotes our method with the coarse global model for final fusion.

Ablation Study of AABB Scale Factor. We present the ablation study of the AABB scale factor in Tab. 3. The scale factor controls the overlap regions across different blocks. It is expected that when we increase the scale factor, the

reconstruction quality can improve since each block can have more observations. However, a too-large scale factor can damage the reconstruction quality. As we can see from the last column in Tab. 3, the PSNR, SSIM, LPIPS when $s_{\text{AABB}} = 1.3$ are worse than those when $s_{\text{AABB}} = 1.2$. This suggests that a moderate scale factor contributes to improved reconstruction quality, while a too-large scale factor may lead to a higher probability of hash collisions and a decrease in performance. We also present the qualitative results of the ablation study on the ‘Caterpillar’ scene in Fig. 7.

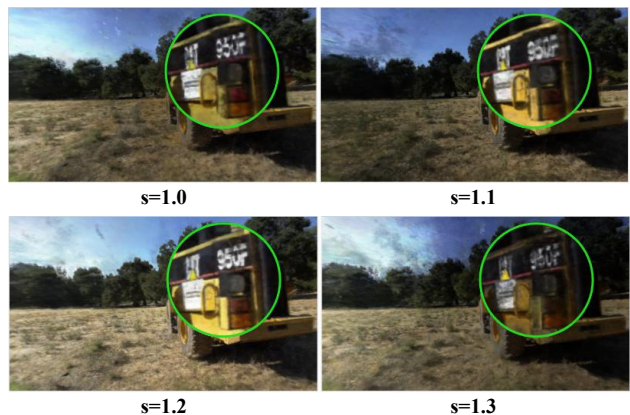


Figure 7. Qualitative results of the **ablation study for the AABB scale factor** on the ‘Caterpillar’ scene of tanks-and-temple dataset.

	$s_{\text{AABB}} = 1.0$	$s_{\text{AABB}} = 1.1$	$s_{\text{AABB}} = 1.2$	$s_{\text{AABB}} = 1.3$
PSNR \uparrow	17.512	17.519	18.146 \bullet	17.564
SSIM \uparrow	0.565 \bullet	0.534	0.555	0.532
LPIPS \downarrow	0.435	0.429	0.406 \bullet	0.428

Table 3. **Ablation studies of the AABB scale factor s_{AABB} .**

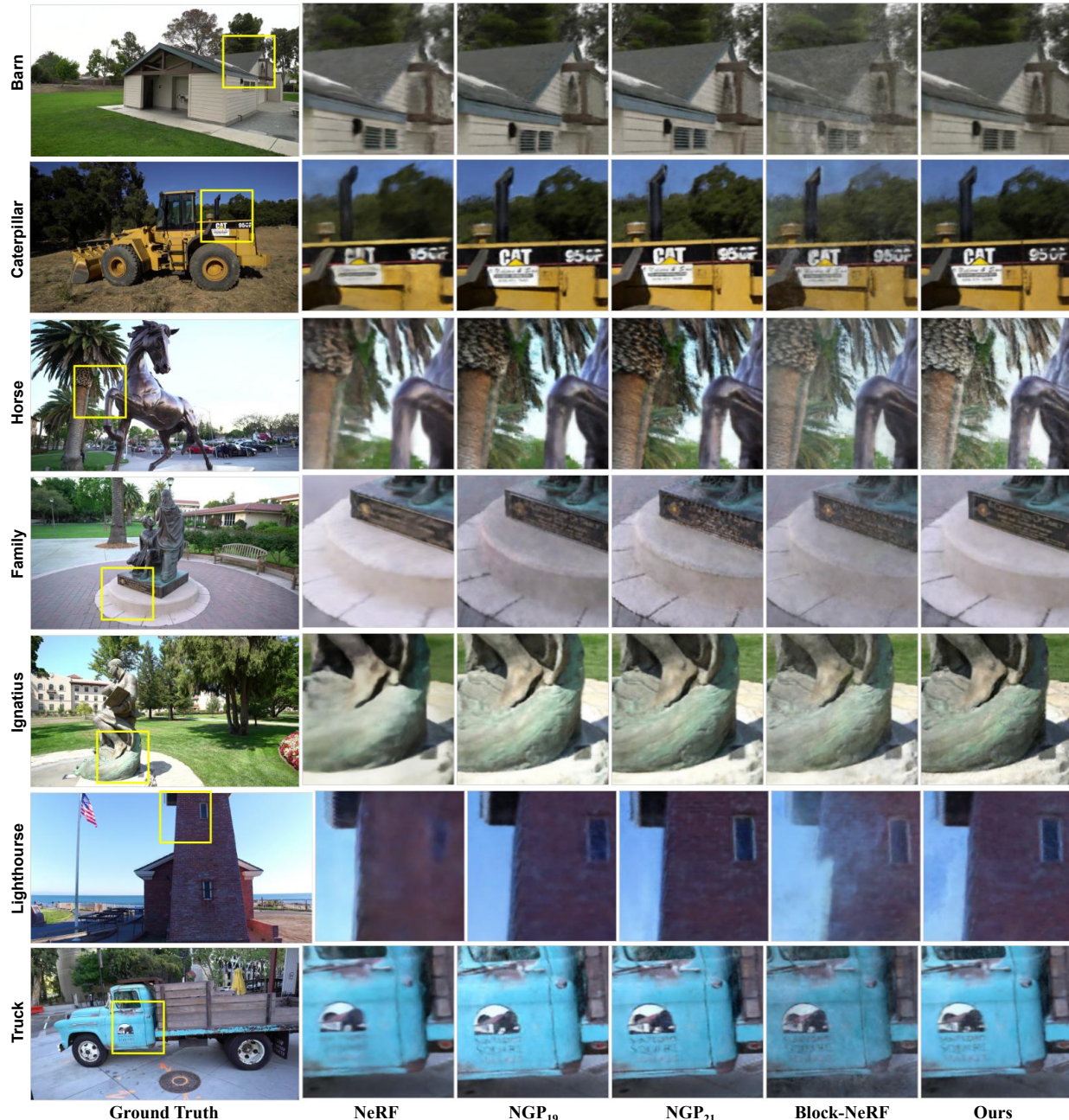


Figure 6. **The qualitative results on tanks-and-templates dataset [22].** From left to right are ground-truth image, vanilla NeRF [32], Instant-NGP [33] with hash table size 2^{19} , InstantNGP [33] with hash table size 2^{21} , Block-NeRF [44], our Scalar-NeRF.

5. Conclusion

We introduced SCALAR-NeRF, a pioneering coarse-to-fine framework designed for the intricate task of large-scale scene reconstruction. Utilizing the encoder-decoder architecture of NeRF, we devised a strategy to decompose expansive scenes into manageable blocks. Our innovation lies in the shared decoder, a component inherited by local blocks from a robust coarse global model and trained comprehensively on all available images. This shared decoding mech-

anism, coupled with the strategic use of the global model in the final fusion step, elegantly addresses challenges posed by imperfect inverse distance weighting functions as observed in approaches like Block-NeRF. The evaluation of our SCALAR-NeRF on the demanding tanks-and-templates dataset reveals a commendable scalability in reconstruction quality. By leveraging the power of a global model and strategically employing it in the fusion stage, we effectively mitigate issues associated with foggy reconstructions.

References

- [1] Relja Arandjelovic, Petr Gronát, Akihiko Torii, Tomás Pajdla, and Josef Sivic. Netvlad: CNN architecture for weakly supervised place recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5297–5307, 2016. 2
- [2] Connelly Barnes, Eli Shechtman, Adam Finkelstein, and Dan B. Goldman. Patchmatch: a randomized correspondence algorithm for structural image editing. *ACM Trans. Graph.*, 28(3):24, 2009. 2
- [3] Jonathan T. Barron, Ben Mildenhall, Matthew Tancik, Peter Hedman, Ricardo Martin-Brualla, and Pratul P. Srinivasan. Mip-nerf: A multiscale representation for anti-aliasing neural radiance fields. In *IEEE/CVF International Conference on Computer Vision*, pages 5835–5844, 2021. 3
- [4] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5460–5469, 2022. 3
- [5] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *CoRR*, abs/2304.06706, 2023. 3
- [6] Anpei Chen, Zexiang Xu, Fuqiang Zhao, Xiaoshuai Zhang, Fanbo Xiang, Jingyi Yu, and Hao Su. Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *2021 IEEE/CVF International Conference on Computer Vision*, pages 14104–14113. IEEE, 2021. 3
- [7] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *Computer Vision - ECCV 2022 - 17th European Conference*, pages 333–350, 2022. 1, 3
- [8] Yu Chen and Gim Hee Lee. Dreg-nerf: Deep registration for neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22703–22713, 2023. 3
- [9] Yu Chen and Gim Hee Lee. Dbarf: Deep bundle-adjusting generalizable neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24–34, 2023. 3
- [10] Yu Chen, Shuhan Shen, Yisong Chen, and Guoping Wang. Graph-based parallel large scale structure from motion. *Pattern Recognition*, 107:107537, 2020. 2, 3, 4
- [11] Yu Chen, Zihao Yu, Shu Song, Tianning Yu, Jianming Li, and Gim Hee Lee. Adasfm: From coarse global to fine incremental adaptive structure from motion. In *IEEE International Conference on Robotics and Automation*, pages 2054–2061. IEEE, 2023. 2
- [12] Kangle Deng, Andrew Liu, Jun-Yan Zhu, and Deva Ramanan. Depth-supervised NeRF: Fewer views and faster training for free. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [13] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinlong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5491–5500, 2022. 2
- [14] Yasutaka Furukawa and Carlos Hernández. Multi-view stereo: A tutorial. *Found. Trends Comput. Graph. Vis.*, 9(1-2):1–148, 2015. 2
- [15] Yasutaka Furukawa, Brian Curless, Steven M. Seitz, and Richard Szeliski. Towards internet-scale multi-view stereo. In *The IEEE Conference on Computer Vision and Pattern Recognition*, pages 1434–1441, 2010. 2
- [16] David Gallup, Jan-Michael Frahm, Philippos Mordohai, Qingxiong Yang, and Marc Pollefeys. Real-time plane-sweeping stereo with multiple sweeping directions. In *IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2007. 2
- [17] Lily Goli, Daniel Rebain, Sara Sabour, Animesh Garg, and Andrea Tagliasacchi. nerf2nerf: Pairwise registration of neural radiance fields. *CoRR*, abs/2211.01600, 2022. 3
- [18] Lily Goli, Cody Reading, Silvia Sellán, Alec Jacobson, and Andrea Tagliasacchi. Bayes’ rays: Uncertainty quantification for neural radiance fields. *CoRR*, abs/2309.03185, 2023. 3, 6, 1
- [19] Mohammad Mahdi Johari, Yann Lepoittevin, and François Fleuret. Geonerf: Generalizing nerf with geometry priors. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18344–18347. IEEE, 2022. 3
- [20] Bernhard Kerbl, Georgios Kopanas, Thomas Leimkühler, and George Drettakis. 3d gaussian splatting for real-time radiance field rendering. *ACM Trans. Graph.*, 42(4):139:1–139:14, 2023. 3
- [21] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations*, 2015. 6
- [22] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics*, 36(4), 2017. 6, 7, 8, 1, 2
- [23] Ruilong Li, Hang Gao, Matthew Tancik, and Angjoo Kanazawa. Nerfacc: Efficient sampling accelerates nerfs. *CoRR*, abs/2305.04966, 2023. 6
- [24] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. BARF: bundle-adjusting neural radiance fields. In *2021 IEEE/CVF International Conference on Computer Vision*, pages 5721–5731. IEEE, 2021. 3
- [25] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. In *Advances in Neural Information Processing Systems 33*, 2020. 2
- [26] Yuan Liu, Sida Peng, Lingjie Liu, Qianqian Wang, Peng Wang, Christian Theobalt, Xiaowei Zhou, and Wenping Wang. Neural rays for occlusion-aware image-based rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7814–7823. IEEE, 2022. 3
- [27] David G. Lowe. Distinctive image features from scale-invariant keypoints. *International Journal Computer Vision*, 60(2):91–110, 2004. 2
- [28] Keyang Luo, Tao Guan, Lili Ju, Haipeng Huang, and Yawei Luo. P-mvsnet: Learning patch-wise matching confidence

- aggregation for multi-view stereo. In *2019 IEEE/CVF International Conference on Computer Vision*, pages 10451–10460, 2019. 2
- [29] Ricardo Martin-Brualla, Noha Radwan, Mehdi S. M. Sajjadi, Jonathan T. Barron, Alexey Dosovitskiy, and Daniel Duckworth. Nerf in the wild: Neural radiance fields for unconstrained photo collections. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7210–7219, 2021. 1, 3, 6
- [30] Quan Meng, Anpei Chen, Haimin Luo, Minye Wu, Hao Su, Lan Xu, Xuming He, and Jingyi Yu. Gnerf: Gan-based neural radiance field without posed camera. In *2021 IEEE/CVF International Conference on Computer Vision*, pages 6331–6341. IEEE, 2021. 3
- [31] Zhenxing Mi and Dan Xu. Switch-nerf: Learning scene decomposition with mixture of experts for large-scale neural radiance fields. In *The Eleventh International Conference on Learning Representations*, 2023. 3
- [32] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *Computer Vision - ECCV 2020 - 16th European Conference*, pages 405–421, 2020. 1, 2, 3, 6, 7, 8
- [33] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM Trans. Graph.*, 41(4):102:1–102:15, 2022. 1, 2, 3, 6, 7, 8
- [34] David Nistér and Henrik Stewénius. Scalable recognition with a vocabulary tree. In *IEEE Computer Conference on Computer Vision and Pattern Recognition*, pages 2161–2168, 2006. 2
- [35] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019*, pages 8024–8035, 2019. 6
- [36] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 10318–10327, 2021. 1
- [37] Konstantinos Rematas, Andrew Liu, Pratul P. Srinivasan, Jonathan T. Barron, Andrea Tagliasacchi, Thomas A. Funkhouser, and Vittorio Ferrari. Urban radiance fields. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12922–12932, 2022. 3
- [38] Barbara Roessle, Jonathan T. Barron, Ben Mildenhall, Pratul P. Srinivasan, and Matthias Nießner. Dense depth priors for neural radiance fields from sparse input views. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 3
- [39] Sara Sabour, Suhani Vora, Daniel Duckworth, Ivan Krasin, David J. Fleet, and Andrea Tagliasacchi. Robustnerf: Ignoring distractors with robust losses. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 20626–20636, 2023. 3
- [40] Johannes L. Schönberger and Jan-Michael Frahm. Structure-from-motion revisited. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pages 4104–4113, 2016. 1, 2
- [41] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Photo tourism: exploring photo collections in 3d. *ACM Trans. Graph.*, 25(3):835–846, 2006. 2
- [42] Mohammed Suhail, Carlos Esteves, Leonid Sigal, and Ameesh Makadia. Light field neural rendering. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8259–8269. IEEE, 2022. 3
- [43] 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. In *Advances in Neural Information Processing Systems*, 2020. 2
- [44] Matthew Tancik, Vincent Casser, Xinchun Yan, Sabeek Pradhan, Ben P. Mildenhall, Pratul P. Srinivasan, Jonathan T. Barron, and Henrik Kretzschmar. Block-nerf: Scalable large scene neural view synthesis. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8238–8248, 2022. 2, 3, 6, 7, 8
- [45] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12912–12921, 2022. 2, 3, 5
- [46] Haithem Turki, Jason Y. Zhang, Francesco Ferroni, and Deva Ramanan. SUDS: scalable urban dynamic scenes. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12375–12385, 2023. 3
- [47] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, Pablo Speciale, and Marc Pollefeys. Patchmatchnet: Learned multi-view patchmatch stereo. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 14194–14203, 2021. 2
- [48] Fangjinhua Wang, Silvano Galliani, Christoph Vogel, and Marc Pollefeys. Itermv: Iterative probability estimation for efficient multi-view stereo. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8596–8605. IEEE, 2022. 2
- [49] Peng Wang, Lingjie Liu, Yuan Liu, Christian Theobalt, Taku Komura, and Wenping Wang. Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021*, pages 27171–27183, 2021. 1
- [50] Qianqian Wang, Zhicheng Wang, Kyle Genova, Pratul P. Srinivasan, Howard Zhou, Jonathan T. Barron, Ricardo Martin-Brualla, Noah Snavely, and Thomas A. Funkhouser. Ibrnet: Learning multi-view image-based rendering. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4690–4699. Computer Vision Foundation / IEEE, 2021. 3

- [51] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.*, 13(4): 600–612, 2004. [6](#)
- [52] Yi Wei, Shaohui Liu, Yongming Rao, Wang Zhao, Jiwen Lu, and Jie Zhou. Nerfingmvs: Guided optimization of neural radiance fields for indoor multi-view stereo. In *2021 IEEE/CVF International Conference on Computer Vision*, pages 5590–5599. IEEE, 2021. [3](#)
- [53] 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. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8296–8306, 2023. [3](#)
- [54] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *Computer Vision - ECCV 2018 - 15th European Conference*, pages 785–801, 2018. [2](#)
- [55] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5525–5534, 2019. [2](#)
- [56] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *2021 IEEE/CVF International Conference on Computer Vision*, pages 5732–5741, 2021. [2](#)
- [57] Richard Zhang, Phillip Isola, Alexei A. Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *2018 IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595. Computer Vision Foundation / IEEE Computer Society, 2018. [6](#)
- [58] Siyu Zhu, Tianwei Shen, Lei Zhou, Runze Zhang, Tian Fang, and Long Quan. Accurate, scalable and parallel structure from motion. *CoRR*, abs/1702.08601, 2017. [2](#)
- [59] Siyu Zhu, Runze Zhang, Lei Zhou, Tianwei Shen, Tian Fang, Ping Tan, and Long Quan. Very large-scale global sfm by distributed motion averaging. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4568–4577, 2018. [2](#), [3](#), [4](#)

SCALAR-NeRF: SCALable LARge-scale Neural Radiance Fields for Scene Reconstruction

Supplementary Material

6. Ablation Study of Coarse Global Model

Our method is evaluated on the large-scale tanks-and-temple dataset. We do not evaluate our method on the city-scale scenes as in Block-NeRF due to GPU memory limitation, since it still requires larger memory for extremely large scenes to produce satisfactory results. The main potential limitation of our method is the reliance on a coarse global model. On larger scenes, the performance of the coarse global model can degrade. To assess the robustness and scalability of our method, we conducted an ablation study on the hash table size of the coarse global model. Given the inherent GPU memory limitations for extremely large scenes, we simulated scenarios where the coarse global model is trained on larger-scale scenes with a reduced hash table size. Specifically, we varied the hash table size of the coarse global model to 17 and 18, while maintaining the hash table size of each local model at 19, as detailed in the main paper.

We present the quantitative results in Table 4 and showcase the performance of our method under different hash table sizes for the coarse global model. Notably, even when the hash table size is reduced, our method consistently outperforms InstantNGP with an equivalent hash table size. More intriguingly, Ours18 surpasses the performance of NGP19, underscoring the potential of our method to scale up to larger scenes. This observation suggests that our approach maintains efficacy and competitiveness, even when confronted with reduced hash table sizes for the global model, signifying its adaptability and promising scalability.

Scenes	NGP ₁₇	NGP ₁₈	NGP ₁₉	Ours ₁₇	Ours ₁₈	Ours ₁₉
Barn	19.033	20.461	20.422	21.733	22.025	21.921
Caterpillar	15.709	17.944	17.220	15.440	18.509	18.280
Courthouse	16.410	15.162	15.088	17.522	15.307	15.990
Family	18.500	17.494	17.606	17.985	17.931	17.838
Horse	17.304	18.879	17.632	17.033	17.376	18.424
Ignatius	15.904	15.956	15.709	16.643	17.029	17.498
Lighthouse	14.934	16.163	15.369	15.030	15.711	16.280
Truck	18.816	17.706	18.562	17.495	17.623	17.984
Average	17.076	17.526 ●	17.201	17.360	17.689 ●	18.027 ●

Table 4. Ablation study of the hash table size on tanks-and-temple [22] dataset. NGP₁₇, NGP₁₈ and NGP₁₉ denote InstantNGP [33] with hash table size 2¹⁷, 2¹⁸ and 2¹⁹, respectively. Ours₁₇, Ours₁₈ and Ours₁₉ denote the hash table size of our coarse global model are 2¹⁷, 2¹⁸ and 2¹⁹, respectively. ● ● and ● respectively denotes the first, second, and third-best results.

7. Uncertainty-based Image Blending

In this section, we also studied the final image blending by leveraging the NeRF uncertainty. The uncertainty performs as a metric that indicates whether the rendered pixel color is reliable or not. Bays' rays [18] compute the uncertainties of the NeRF model after training to remove floaters. The intuition behind the method is that by adding a small perturbation to the point coordinate \mathbf{x} , the underlying geometry of the learned model should not change too much. They therefore adopt a deformation field to do this:

$$\mathcal{D}_{\theta(\mathbf{x})} = \text{Trilinear}(\mathbf{x}, \theta). \quad (11)$$

The volume density and per-pixel color are therefore updated as:

$$\sigma, \mathbf{c} = \text{Decoder}(\text{Encoder}(\mathbf{x} + \mathcal{D}_{\theta(\mathbf{x})})). \quad (12)$$

The uncertainty is then computed by the inverse of the diagonal components of the approximated hessian:

$$\Sigma \approx \text{diag}\left(\frac{2}{R} \sum_{\mathbf{r}} \mathbf{J}_{\theta}(\mathbf{r})^{\top} \mathbf{J}_{\theta}(\mathbf{r}) + 2\lambda \mathbf{I}\right)^{-1}, \quad (13)$$

where R is the total number of rays \mathbf{r} , \mathbf{J}_{θ} is the Jacobian matrix of the rendered color w.r.t. the model parameters θ .

Once we computed the uncertainty for the rendered pixel, we can naturally pick out the one that has the least uncertainty for image blending. Though Bayes' rays has shown promising results in removing floaters, we find the uncertainties are not reliable enough for high-fidelity image blending. As we can see in Fig. 8, the uncertainty-based method can also give very low uncertainty for areas that are under-constrained. Moreover, for ambiguous areas, such as the sky, the uncertainty can be always high, which makes it difficult to filter wrong results from different blocks. The quantitative results provided on Tab. 5 also validate the unreliability of the uncertainty-based method.

Scenes	PSNR ↑		SSIM ↑		LPIPS ↓	
	Bayes' Rays	Ours	Bayes' Rays	Ours	Bayes' Rays	Ours
Barn	14.517	21.921	0.451	0.639	0.577	0.378
Caterpillar	14.209	18.280	0.359	0.541	0.537	0.420
Courthouse	12.226	15.990	0.419	0.500	0.630	0.534
Family	13.517	17.838	0.304	0.471	0.595	0.380
Horse	13.875	18.424	0.410	0.648	0.574	0.338
Ignatius	12.811	17.498	0.235	0.458	0.548	0.405
Lighthouse	11.784	16.280	0.367	0.548	0.678	0.537
Truck	13.223	17.984	0.353	0.527	0.543	0.387
Average	13.270	18.027	0.362	0.542	0.585	0.422

Table 5. Quantitative comparisons of uncertainty-based method and ours on tanks-and-temple [22] dataset.

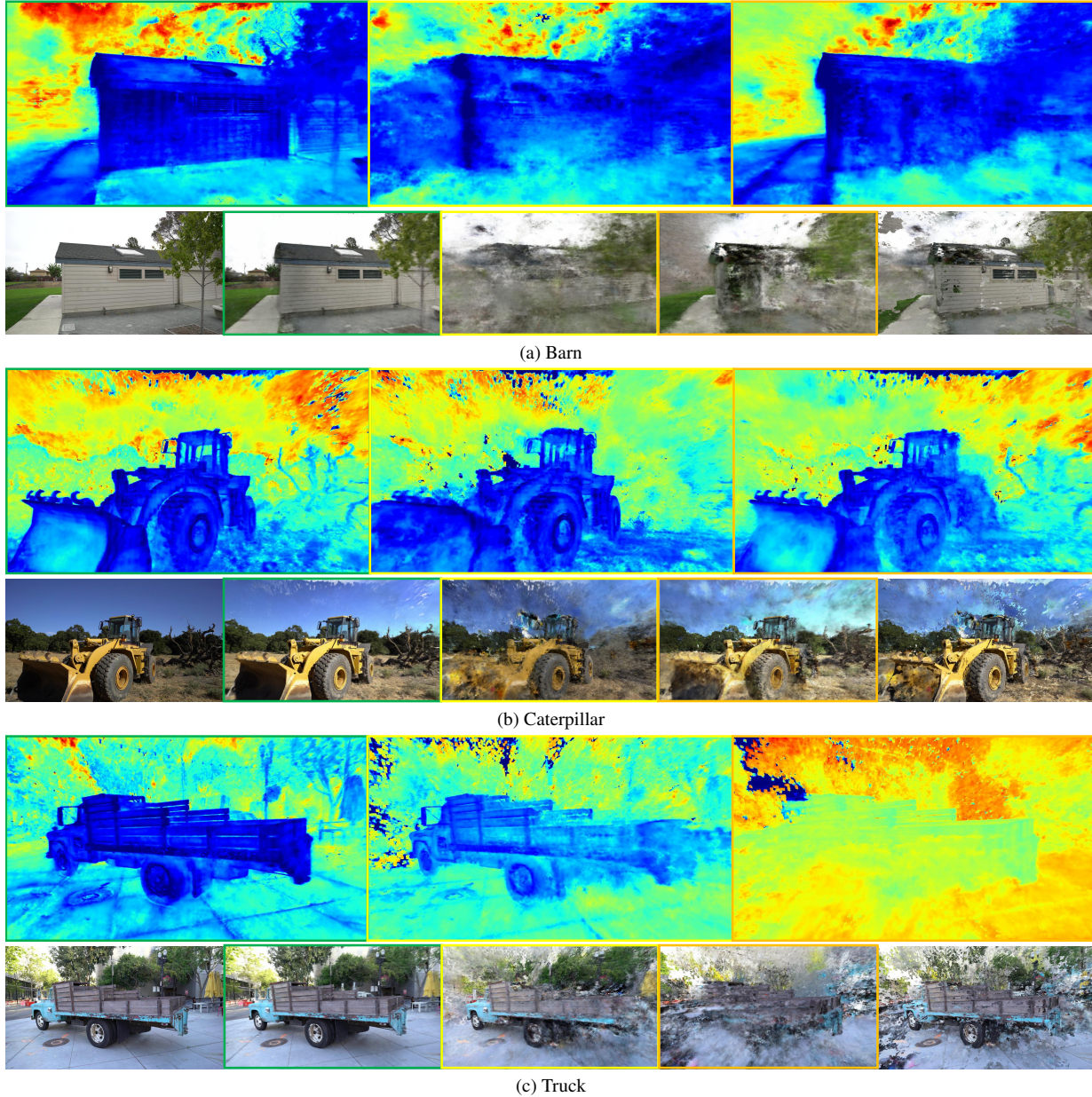


Figure 8. **The qualitative results of uncertainty-based image blending on tanks-and-temples dataset [22].** For each sub-figure, the first row represents the uncertainty heatmap for each block, the second row represents the ground-truth image, the images rendered from each block, and the final fused images.

8. Limitations and Future Work

Our approach excels in producing higher-quality images compared to Block-NeRF. However, it comes with a trade-off, as it relies on initializing a coarse NeRF model. Consequently, the training time for our method is marginally longer than that of Block-NeRF, especially when employing the Mip-NeRF backbone with InstantNGP. Additionally, during inference, our method necessitates the gener-

ation of a reference image from the coarse model, introducing a slight slowdown in inference speed.

While our work aims to strike a balance between computational efficiency and reconstruction quality, addressing these limitations, particularly in the context of larger scenes, and reducing the dependence on coarse global models presents avenues for future exploration. These considerations underscore the necessity for ongoing refinement and adaptation to diverse environmental scales and conditions.

9. Social Impact

Our method presents a significant stride forward in the field of large-scale scene reconstruction. Beyond its technical prowess, the potential social impact of this innovation is multifaceted.

- **Enhanced Visualizations.** The ability to reconstruct large-scale scenes with improved fidelity and scalability can have a profound impact on various industries. From urban planning to environmental monitoring, SCALAR-NeRF's capacity to generate high-quality visualizations can aid decision-makers and stakeholders in better understanding and analyzing complex spatial data.
- **Accessibility and Resource Efficiency.** The scalability of SCALAR-NeRF, demonstrated by its capability to reconstruct extensive scenes on a single GPU, holds the promise of democratizing access to advanced reconstruction technologies. This accessibility could empower researchers, educators, and professionals who may have been limited by resource constraints, fostering innovation and collaboration.
- **Advancements in Cultural Preservation.** Large-scale scene reconstruction is crucial in fields like archaeology and cultural heritage preservation. SCALAR-NeRF's ability to faithfully capture detailed scenes can contribute to the preservation and documentation of historical sites and artifacts, aiding in cultural conservation efforts globally.
- **Improved Urban Planning.** As urban areas expand, the demand for sophisticated tools in urban planning intensifies. SCALAR-NeRF's detailed reconstructions can provide urban planners with realistic and data-rich representations of large-scale environments, facilitating more informed decision-making in areas such as infrastructure development and disaster preparedness.
- **Training and Simulation.** In domains like autonomous vehicles and robotics, realistic scene reconstructions are invaluable for training and simulation. SCALAR-NeRF's ability to efficiently handle large-scale scenes can contribute to advancements in autonomous systems by providing realistic training environments.

While the immediate impact is within the realms of research and technology, the broader social implications lie in the democratization of advanced reconstruction capabilities, fostering innovation, and contributing to the understanding and preservation of our physical environment and cultural heritage.