

C^3 -NeRF: Modeling Multiple Scenes via Conditional-cum-Continual Neural Radiance Fields

Prajwal Singh, Ashish Tiwari, Gautam Vashishtha & Shanmuganathan Raman

{singh_prajwal, ashish.tiwari, gautam.pv, shanmuga}@iitgn.ac.in

CVIG Lab, IIT Gandhinagar, Gujarat, India

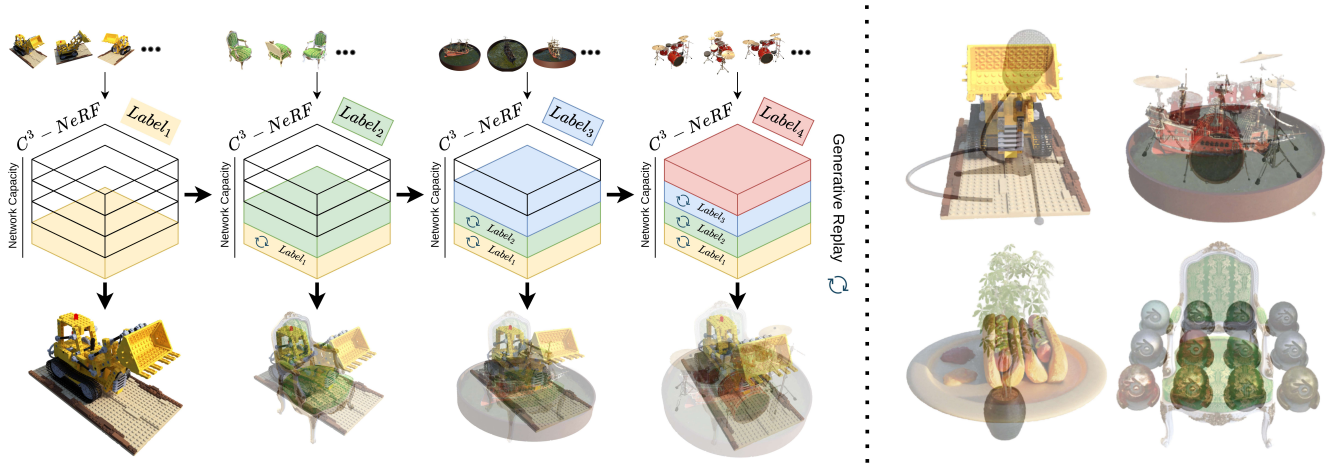


Figure 1: C^3 -NeRF Framework: (Left) The figure illustrates the proposed Conditional-cum-Continual learning framework for modeling multiple scenes in a single neural radiance field. Each new scene is assigned a pseudo label and trained continually while preserving the information of previously learned scenes via generative replay. At any point in time, a specific learned scene(s) can be rendered by conditioning C^3 -NeRF on the associated pseudo-label(s). (Right) The figure visualizes the learned multi-scene neural radiance field. All the eight scenes of the NeRF Synthetic 360° [MST*21] dataset modeled by C^3 -NeRF are shown as separate pairs to avoid clutter.

Abstract

Neural radiance fields (NeRF) have exhibited highly photorealistic rendering of novel views through per-scene optimization over a single 3D scene. With the growing popularity of NeRF and its variants, they have become ubiquitous and have been identified as efficient 3D resources. However, they are still far from being scalable since a separate model needs to be stored for each scene, and the training time increases linearly with every newly added scene. Surprisingly, the idea of encoding multiple 3D scenes into a single NeRF model is heavily under-explored. In this work, we propose a novel conditional-cum-continual framework, called C^3 -NeRF, to accommodate multiple scenes into the parameters of a single neural radiance field. Unlike conventional approaches that leverage feature extractors and pre-trained priors for scene conditioning, we use simple pseudo-scene labels to model multiple scenes in NeRF. Interestingly, we observe the framework is also inherently continual (via generative replay) with minimal, if not no, forgetting of the previously learned scenes. Consequently, the proposed framework adapts to multiple new scenes without necessarily accessing the old data. Through extensive qualitative and quantitative evaluation using synthetic and real datasets, we demonstrate the inherent capacity of the NeRF model to accommodate multiple scenes with high-quality novel-view renderings without adding additional parameters. We provide implementation details and dynamic visualizations of our results in the supplementary file.

CCS Concepts

• Computing methodologies → Rendering;

1. Introduction

Neural Radiance Fields (NeRF) generate novel photorealistic views of a scene from a sparse set of input images by implicitly modeling the 3D scene via Multi-Layer Perceptrons (MLPs) and dif-

ferentiable volumetric rendering. Despite its strengths, NeRF and its variants typically struggle with extensive per-scene optimizations and prolonged training time. Several recent advancements have been made to improve different limiting aspects of NeRF,

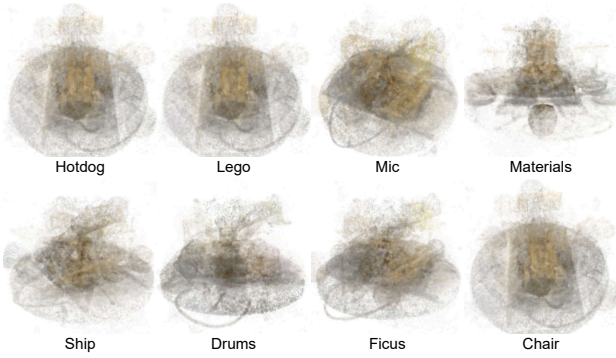


Figure 2: Vanilla Instant-NGP [MESK22]. We combine all the different scenes of the NeRF Synthetic 360° [MST*21] dataset and train Instant-NGP. The results show that neural hash encoding of scene coordinates may not be sufficient to preserve information across multiple scenes.

i.e., expedite NeRF’s training process using caching or employing sparse voxel grids [WLL*21, FKYT*22, YYTK21], reduce inference times [CFHT23, WRB*23], model compression [GCML23, CXG*22, TCWZ22, SP24, FXW*23], scene editing [WWQQ23, BZY*23], generalization [HLX*23, YYTK21, WWG*21, ZBS*22, YHL*23, CYM*23], and few-shot learning [LCK*24, YPW23, WCLL23], to list a few. However, the ability of NeRF models to encode multiple scenes with the same set of parameters has remained heavily under-explored. Furthermore, over the recent years, NeRFs have become astonishingly valuable 3D resources. However, managing them is still tedious since each scene needs a separate model to be trained from scratch, linearly increasing the storage and training time for every new scene. Representing multiple scenes using the same number of parameters as for a single-scene neural radiance field with extra room to accommodate new unseen scenes would offer better scalability and applicability. Furthermore, there is no need to store the training data as it can be generated at point of time via scene conditioning and generative replay.

Resourceful modeling of multiple scenes requires certain fundamental properties - (a) *photo-realistically render all the observed scenes at any point of time*, (b) *efficiently accommodate new scenes without the need for training data of the previous scenes*, (c) *without effecting the learned knowledge of the previous scenes*, and (d) *without increasing the number of parameters*.

We derive inspiration from the benefits of conditioning and continual learning to satisfy these requirements. Conditioning allows for selecting the appropriate parameter subspace for each scene such that, given a condition label, one can render the specific scene from the trained NeRF. While continual learning methods avoid catastrophic forgetting, they cannot be directly integrated with NeRF to reconstruct multiple scenes due to the underlying differences. To better understand this, let us consider image classification models, for instance, where the same image seldom has different labels across tasks. However, in NeRF, the same position and view direction across different scenes (tasks) will have different densities and colors. Moreover, neural networks capture

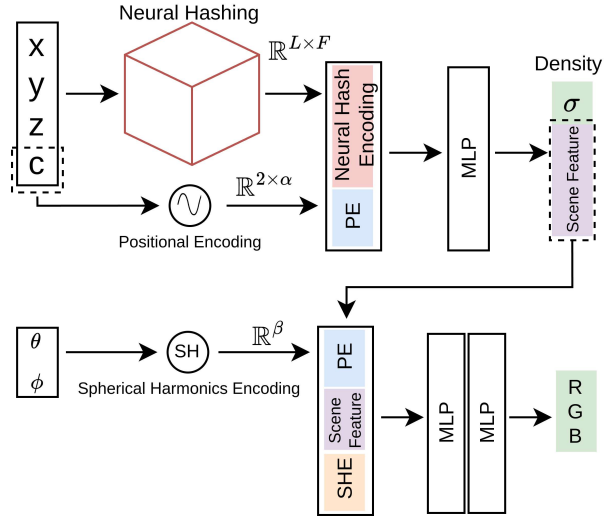


Figure 3: Multi-scene Architecture. The proposed modification to the Instant NGP architecture [MESK22] incorporates information from multiple scenes using pseudo labels. These pseudo labels are encoded with neural hashing and scene coordinates, followed by positional encoding. The encoded labels are then concatenated with the neural hash representation to enhance performance, as discussed in the ablation study (Section 4.10).

generic, reusable features in classification, making continual learning easier. In contrast, NeRF models generally focus on scene-specific details and lack feature re-usability across different scenes. Interestingly, to use the benefits of feature reuse, some methods [CXZ*21, HLX*23, WWG*21, BDH*23, YHL*23] have deployed convolutional feature extractors that learn local and global scene features [PFS*19, YYTK21, BMV*22], attention modules, 3D cost volume features, and other advanced neural networks with NeRF focusing on either faster optimization and/or better generalizability. From another perspective, these additional neural modules provide neural scene conditioning, i.e., the convolutional features of images of the scene help NeRF identify the underlying scene uniquely. However, they do not specifically model multiple scenes in a single NeRF and contain additional parameters in these auxiliary feature extractors. Moreover, these methods are not suited for accommodating new scenes incrementally and continue to suffer from catastrophic forgetting of the previously trained scenes. Although there have been certain attempts to combine continual learning with NeRF [CLBL22, CM23, PDBW23, ZLCX23], they apply only to single-scene reconstruction with multiple appearance or geometry changes over time. Unlike these methods, we aim to mine the potential of NeRF models to represent multiple scenes with the same set of parameters by combining the benefits of conditioning and continual learning.

In this work, we propose C^3 -NeRF - a Conditional-cum-Continual paradigm to model multiple scenes in a single Neural Radiance Field. The design of C^3 -NeRF strategically harnesses the learning ability and training speed of Instant-NGP [MESK22]. Specifically, C^3 -NeRF embeds scene conditioning directly into



Figure 4: Qualitative demonstration of the quality of rendered images across different views through C^3 -NeRF over the scenes from the NeRF Synthetic 360° dataset [MST*21]. Best viewed in pdf with zoom.

the hashing mechanism of Instant-NGP [MESK22] using pseudo-labels to maximally utilize a single NeRF’s parameter space to accommodate multiple scenes. Through simple scene conditioning, our method provides a pioneering ability to allow for flexibility to perform multi-scene optimization, where the same set of model parameters can be made to effectively learn new scenes while maintaining high fidelity in rendering previously encountered ones. Our approach provides a multi-scene optimization without needing additional pre-trained priors or feature extraction blocks. Moreover,

learning under the generative-replay paradigm, C^3 -NeRF does not require the availability of the training data of the previously encountered scene and preserves their rendering fidelity to a large extent while accommodating new scenes. Thus, C^3 -NeRF is more feasible for practical applications, mainly where training data and resources may be limited, to model multiple 3D scenes without extra computational overhead or extensive retraining.

Contributions. The key features and contributions of our work are as follows.

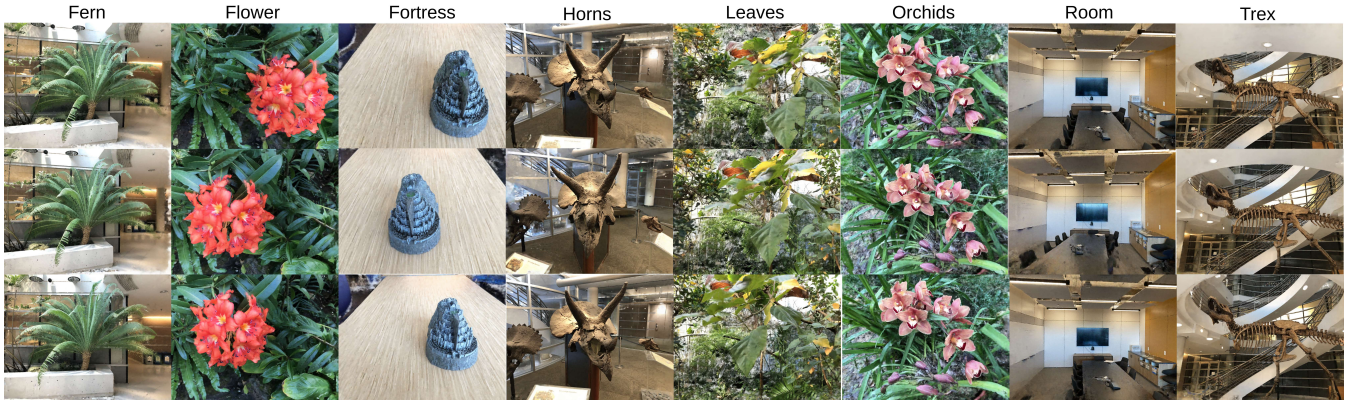


Figure 5: *Qualitative demonstration of the quality of rendered images across different views through C^3 -NeRF over the scenes from the Real Forward-Facing dataset [MSOC*19]. Best viewed in pdf with zoom.*

- C^3 -NeRF takes the first steps toward encoding multiple scenes in a single conditional-cum-continual neural radiance field that adapts to new scenes over time with minimal, if not no, loss in rendering quality in previously encountered scenes. Furthermore, generative replay in C^3 -NeRF removes the dependence on the training data of the previously encountered scenes.
- We achieve rapid adaptation across diverse scenes via a simple and effective scene conditioning through the hashing mechanism of Instant-NGP [MESK22] that naturally offers high rendering fidelity and enhanced training speed than conventional NeRF.
- C^3 -NeRF does not depend on the need for any additional pre-trained priors or feature extractors, avoiding any additional overhead for modeling multiple scenes.
- Through comprehensive experimental evaluations, we demonstrate that C^3 -NeRF consistently outperforms existing related methods over multi-scene rendering quality, training efficiency, and adaptability to new scenes.

Interestingly, we observe that even the vanilla NeRF, under the proposed conditional-cum-continual paradigm, can model multiple scenes, except for the fact that vanilla NeRF takes a huge amount of training time and offers lower performance than the proposed C^3 -NeRF.

Organization. We start by exhaustively discussing the existing related literature in Section 2. We highlight the key shortcomings and underlying differences in the objectives and experimental settings compared to ours. We then describe the basic premise of learning neural radiance fields in Section 3.1 and discuss the proposed framework in Section 3. We perform exhaustive experiments and provide detailed qualitative and quantitative results in Section 4. Finally, we discuss our work’s key features, limitations, and future scope in Section 5.

2. Related Work

Novel View Synthesis. Synthesizing novel views of a scene to recreate unseen perspectives from a limited set of 2D observations is of broad interest in computer vision and computer graphics [TTM*22]. Traditional approaches without explicit 3D struc-

ture representation require densely sampled viewpoints to accurately render a novel scene with photo-realism [LH23, DTM23, GGSC23]. A few methods have used voxel-based representation to encode a 3D spatial structure. However, these methods suffer from restricted spatial resolution and substantial memory demands [BLRW16, LDG18, SG18, WWX*17]. To alleviate the memory overhead, some methods use compact implicit representations such as signed distance functions (SDF) that map spatial coordinates to 3D signed distance to the surface [MON*19, PFS*19, NMOG20].

Recently, Neural Radiance Fields (NeRFs) [MST*21] have emerged popular in implicitly encoding scene geometry and novel view synthesis. NeRF uses multi-layer perceptrons (MLPs) to estimate color and density at a 3D scene point and allows photorealistic view synthesis via differentiable volumetric rendering. However, they tend to suffer from long per-scene optimization times. To reduce training times, methods like Plenoxels [FKYT*22] and Instant-NGP [MESK22] have introduced more efficient training mechanisms by leveraging spherical harmonics, multi-resolution hash-encoding, and fully fused MLPs. Furthermore, several other variants and extensions of NeRFs have been proposed with objectives such as faster inference [CFHT23, WRB*23], 3D scene editing [WWQQ23, BZY*23], using a sparse set of images [LCK*24, YPW23, WCLL23], etc. Despite such advancements, these approaches continue to focus on single-scene optimization and do not consider handling multiple scenes in their parameter space. Another stream of works [HLX*23, YYTK21, WWG*21, ZBS*22, YHL*23, CYM*23] has focused on generalization by learning a single NeRF model on the fly by pre-training on a set of multiple different scenes. However, they still require fine-tuning on unseen scenes to produce high-quality renderings and suffer from severe forgetting of previously learned scenes.

Moreover, these approaches often depend on CNN-based or 3D cost volume-based feature extractors. While the features learned through such feature extractors can be used as scene conditioners, they tend to increase the computational and memory overhead due to the increased number of learnable parameters. In contrast, we propose representing multiple scenes using the same set of NeRF parameters under a conditional-cum-continual paradigm. Specifi-



Figure 6: Qualitative demonstration of the quality of rendered images across different views through C^3 -NeRF over the scenes from the Tanks and Temples dataset [KPZK17]. It is best viewed in PDF with Zoom.

cally, we apply scene conditioning via simple pseudo labels while training NeRF. These labels eventually help re-render the learned scenes with high fidelity at a specific later time without additional computational overhead.

NeRF and Continual Learning. Continual learning presents a critical challenge in machine learning, characterized by the need to learn from sequential datasets without forgetting previous knowledge [Rob95]. In the context of NeRFs, traditional continual learning strategies have centered around managing distribution shifts without retaining historical data. Existing methodologies are broadly categorized into parameter isolation, regularization, and replay [LH17, ML18, KPR*17, RABT17, SLKK17, CM23]. Parameter isolation methods allocate distinct neurons for new tasks, preserving old information but at the cost of increased network size or limited learning capacity due to the availability of limited neurons for new tasks. Regularization strategies penalize changes to critical parameters, which becomes challenging to identify and scale across multiple domains. Replay approaches, including generative models, reintroduce historical data during training to mitigate forgetting [SLKK17, CM23]. However, as discussed earlier, these methods cannot be directly integrated with NeRF for multi-scene modeling due to the vast domain gap between different scenes. Hence, the parameters cannot be reused efficiently with these methods. Although these techniques are effective for within-scene transitions [PDBW23, ZLCX23, CM23], they fail in multi-scene continual learning. Thus, a few initial works integrate continual learning with NeRF for tasks such as localization and mapping (SLAM) using replay [SLOD21, CCW*23, DSQ*24].

Moreover, they assume that the previous training data is always

available to solve optimization problems over all keyframes for bundle adjustment. Other methods, such as MEIL-NeRF [CLBL22] and CLNeRF [CM23], reconstruct the scene with a sequence of multiple partial scans. However, they are focused on a single scene and handle forgetting only partially learned regions of a scene. Unlike these methods, we cater to continually accommodating multiple scenes. Concurrent to our work is [WWW*24], which attempts to continually encode multiple scenes in NeRF using a hyper-network to estimate weights of the MLP in NeRF. Although they have a similar goal of learning multiple scenes with NeRF, they represent scenes as a linear combination of a cross-scene weight matrix and a set of scene-specific weight matrices generated from a global parameter generator, again relying on separate auxiliary networks. Moreover, the number of parameters increases with every newly added scene and slow training, unlike ours, where the number of parameters is fixed and with fast convergence.

It is important to note that although our objective will have an upper bound on the number of scenes that can be modeled using a fixed set of parameters without forgetting, we found that this number is large enough for real-world applications. Finding such an upper bound is not part of the scope of this work, and we wish to explore that in the near future. One can also think of reducing the required number of parameters to model a single scene. Interestingly, there is a series of work that explores model compression of NeRF using weight quantization [GCML23], binarization [SP24], low-rank approximation [CXG*22, TCWZ22], or knowledge distillation [FXW*23]. However, these works' objectives differ from ours, with the key idea being maximally utilizing the given parameter space to encode as many scenes as possible with high fidelity.

3. Method

In this section, we will discuss the proposed framework C³-NeRF for modeling multiple scenes in a conditional-cum-continual paradigm over the set of parameters that thus far have been used to model only single scenes in a neural radiance field. C³-NeRF harnesses the capability of Instant-NGP [MESK22] that offers faster training speed - an essential requirement in conditional-cum-continual learning of scenes. We start with a brief overview of learning neural radiance fields and multi-resolution hash encoding (MHE) proposed in Instant-NGP. We then introduce and explore the conditional capability of Instant-NGP and describe the integration of generative replay [SLKK17] to conditional Instant-NGP to extend it to a continual learning paradigm and enable progressive learning of new scenes without forgetting the information of the previously encountered scenes. We shall also see that generative replay helps inherently mitigate the performance reduction due to the difference in camera parameters of the scenes. Finally, we discuss the loss functions used to train the complete framework. Interestingly, the proposed conditional-cum-continual learning produces equally good, and sometimes better, results than baselines that handle only a single scene at a time with minimal modification to the Instant-NGP [MESK22].

3.1. Background

3.1.1. Neural Radiance Fields

Mildenhall *et al.* [MST*21] proposed a method for synthesizing unseen multi-view images from a sparse set of images. The framework consists of a multi-layer perceptron (MLP), which takes a 5D input vectors - 3D position vector of a scene point (x, y, z) and the associated viewing direction (θ, ϕ) , and outputs the density (σ) and color (c) at that point. These 3D points are sampled along the rays projected from each pixel in the image and the camera center \mathbf{o} , such that $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$. Here, t is the sampled distance along the ray, and \mathbf{d} is the direction vector from the camera center to the pixel in the image. To render the color along the sampled ray, a discrete version of volumetric rendering is used [Max95], as described in Equation 1

$$\hat{C}(r) = \sum_{i=1}^N w_i c_i \quad (1)$$

Here, $w_i = T_i \alpha_i$ is the weight for each of the N points sampled on the ray and accumulates the transmittance T_i and alpha value $\alpha_i = (1 - \exp(-\sigma_i \delta_i))$. Here, $\delta_i = t_{i+1} - t_i$ is the distance between adjacent samples, and α values are similar to those used in alpha compositing. Broadly speaking, the transmittance value accounts for the light attenuated or transmitted, thus far, i.e., $T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \delta_j\right)$.

To train the MLP network, the squared error between the rendered pixel $\hat{C}(r)$ and ground-truth pixel $C(r)$ is calculated, as per Equation 2.

$$\mathcal{L} = \sum_{r \in \mathcal{R}} \|\hat{C}(r) - C(r)\|_2^2 \quad (2)$$

Here, \mathcal{R} is a set of rays passing through each pixel in the im-

Algorithm 1: Learning new scenes using Conditional-cum-Continual NeRF (C³-NeRF)

Require: I_{new} : list of new scene images
Require: l_{new} : pseudo label for new scene
Require: L_{prev} : list of pseudo labels for previous scenes
Require: F_θ : NeRF network
Require: $C_{I_{new}}$: camera parameters of new scene
Require: k : number of images to render for previous scenes
 /* Render previous scenes and store in I_{prev} */
 Initialize $I_{prev} \leftarrow \{\}$
for $l_{prev}^i \in L_{prev}$ **do**
 $render_image \leftarrow F_\theta(l_{prev}^i, C_{I_{new}}, k)$
 $I_{prev}.append(render_image)$
 $T_{images} \leftarrow concat(I_{new}, I_{prev})$ // train images list
 $T_{labels} \leftarrow concat(L_{new}, L_{prev})$ // train labels list
 /* TrainNeRF() train the F_θ on given scene list
 and labels */
 $F_\theta^{updated} \leftarrow TrainNeRF(F_\theta, T_{images}, T_{labels}, C_{I_{new}})$
Output: Updated $F_\theta^{updated}$ on new scene

age and the camera center. The sinusoidal positional encoding is applied over the 5D input vector before passing it to the MLP.

3.1.2. Multi-Resolution Hash Encoding

Müller *et al.* [MESK22] introduced the concept of multi-resolution hash encoding (MHE) to enhance both the reconstruction accuracy and training efficiency of neural radiance fields (NeRFs) while maintaining minimal computational overhead. Unlike the traditional positional encoding used in NeRF, MHE employs a trainable multi-level 3D grid structure. The position of a sampled point is encoded by interpolating the features (\mathbb{R}^F) located at the vertices of the grid cell containing the point. These interpolated features are then fed into a multi-layer perceptron (MLP) network, which subsequently predicts both the density (σ) and the color (c) of the point.

The grid is organized into L levels, each containing a hash table of size $T \times F$, where T represents the number of features and F is the dimensionality of each feature. To efficiently map the grid coordinates, a spatial hash function $h(x)$, based on the approach of [THM*03], is utilized:

$$h(x) = \left(\bigoplus_{i=1}^d x_i \pi_i \right) \bmod T \quad (3)$$

Here, \oplus denotes the bitwise XOR operation, and π_i are distinct large prime numbers.

3.2. Conditional NeRF

To condition the Instant-NGP network [MESK22], we utilize pseudo integer labels or classes such as $\{1, 2, 3, \dots\}$ - to differentiate each scene. Figure 3 provides an overview of the conditional Instant-NGP pipeline. In this approach, we first obtain Neural Hash

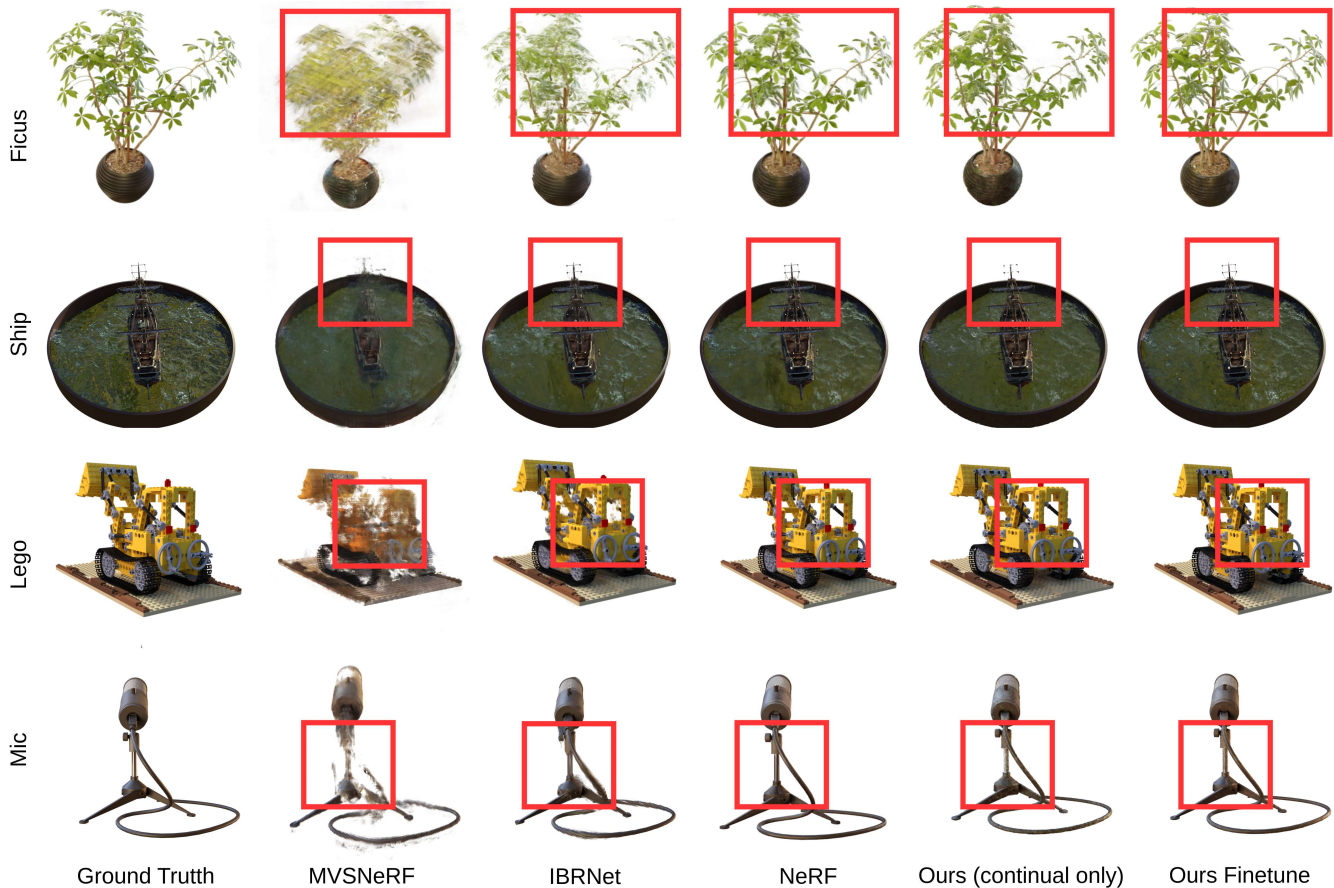


Figure 7: Qualitative Comparison on the Nerf Synthetic 360° dataset [MST*21]. C^3 -NeRF is compared with MVSNeRF [CXZ*21] and IBRNet [WWG*21] (both with per-scene fine-tuning), and vanilla NeRF [MST*21] with per-scene optimization. C^3 -NeRF is observed to model scenes with minimal artifacts compared to those present in the other methods. The quality is enhanced further upon additional per-scene fine-tuning over C^3 -NeRF. The marked regions highlight the differences and are best viewed in PDF with Zoom.

Encoding (NHE) by combining the scene coordinates (x, y, z) and the pseudo label through neural hashing. Apart from neural hashing, the pseudo labels are also processed with sinusoidal positional encoding (PE) [TSM*20] $\psi \in \mathbb{R}^{2\alpha}$, where $\alpha = 2$ denotes the number of frequencies used in our work. Neural hash encoding and the positionally encoded pseudo label are then concatenated and fed into a single-layer MLP to predict density (σ) and scene feature.

Neural hashing when conditioned on the pseudo label, plays a crucial role in segregating the scenes within the shared representation space. As shown in Figure 2, we tested this approach by training Instant-NGP [MESK22] on a mixed dataset of scenes without any labels. The apparent intermixing of scene information in the experiment validated our hypothesis that the neural hashing of just the scene coordinates does not suffice to distinguish between scenes in a multi-scene setting. To achieve unique scene representations, additional conditioning is required.

The predicted scene feature is then passed to the next stage of the Instant-NGP pipeline to compute the color for each input coordinate (x, y, z) . Specifically, the scene feature is concatenated

with the encoded viewing direction and the encoded pseudo label to serve as input into a two-layer MLP to predict the final color values. We encode the viewing direction via spherical harmonics encoding (SHE) [RH01], further refining the rendering process.

3.3. Conditional-cum-Continual NeRF

Simply conditioning the NeRF model as proposed above has certain limitations related to differences in camera parameters, as described in Section 4. Therefore, we deploy the continual learning paradigm and adapt conditional NeRF to a new scene via generative replay [SLKK17]. The learning paradigm of C^3 -NeRF is as follows. We encode a scene into conditional Instant-NGP with the help of a pseudo label. Before encoding the new scene, we render all the previously observed scenes with the camera parameters of the new scene, add them to the training images of the new scene, and update the parameters with this new training set. This strategy helps overcome the dependency on the previous training datasets and models each observed scene with high fidelity. The algorithm 1 shows the proposed training strategy for C^3 -NeRF.

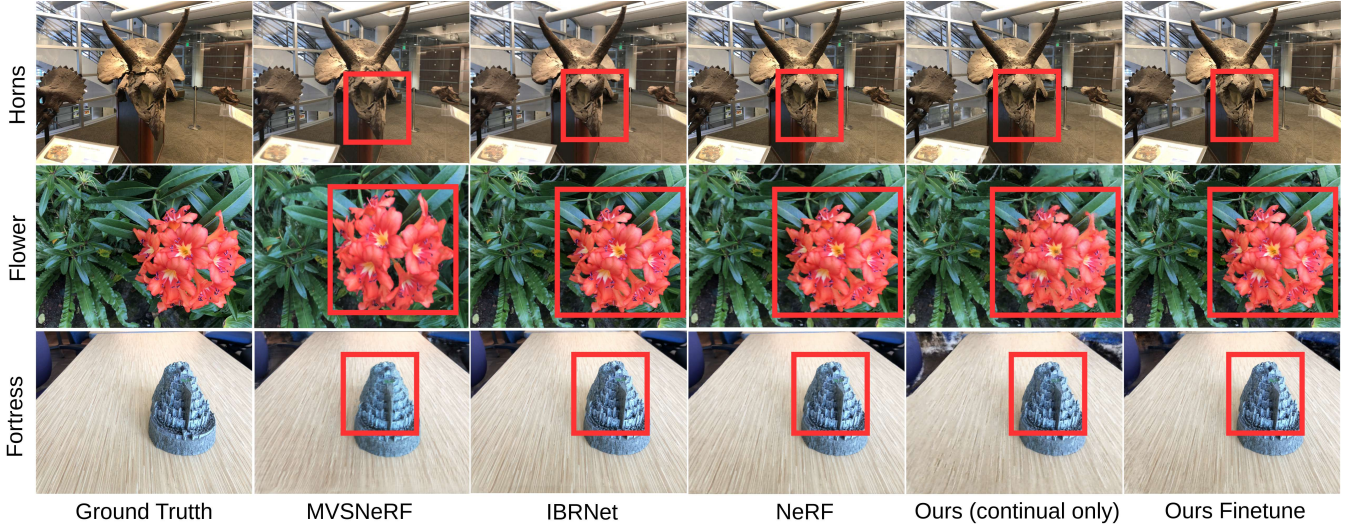


Figure 8: *Qualitative Comparison on the Real Forward Facing (LLFF) dataset [MSOC*19]. We compare C^3 -NeRF with MVSNeRF [CXZ*21] and IBNet [WWG*21] (both with per-scene fine-tuning), and vanilla NeRF [MST*21] with per-scene. C^3 -NeRF performs nearly equivalent to vanilla NeRF and performs better upon per-scene fine-tuning. The marked regions show the differences and are best viewed in PDF with Zoom.*

In Figure 9, we present the t-SNE map [VdMH08] to visualize the learned multi-scene representation space of C^3 -NeRF across eight scenes from the NeRF Synthetic 360° dataset [MST*21]. The overlapping scene clusters suggest that the parameter space is shared across multiple scenes. This indicates that a subset of parameters can encode shared information across different scenes rather than being dedicated solely to a specific one. Specifically, we extract the neural hash encoding (as described in Figure 3) from one image per scene, each rendered from the same view. Since the image consists of 800×800 pixels and a lot of the pixels would correspond to the background, we divide the image into 50 patches, resulting in 16 patches per image, and extract the MLP features from these patches instead of every pixel. These features across all the 8 scenes are visualized via t-SNE map [VdMH08].

3.4. Loss and Regularization

Floater represents a significant challenge in NeRF implementations, typically manifesting as disconnected, dense spatial regions near the camera plane. We have implemented two essential regularization techniques to address this issue and ensure stable C^3 -NeRF training.

Distortion: Following the work of [BMV*22, WT23], we used distortion loss for compact point distribution on the camera ray:

$$\mathcal{L}_{dist} = \frac{1}{d(r)} \left(\sum_i w_i w_j \left| \frac{t_i + t_{i+1}}{2} - \frac{t_j + t_{j+1}}{2} \right| + \frac{1}{3} \sum_{i=1}^N w_i^2 (t_{i+1} - t_i) \right) \quad (4)$$

where, $d(r) = \frac{\sum_{i=1}^N w_i t_i}{\sum_{i=1}^N w_i}$ is depth along each ray. The first part of \mathcal{L}_{dist} minimizes the weighted distance between all pairs of interval

midpoints. The second part focuses on minimizing the weighted size of each interval. Jointly, the weights on the ray are encouraged to be compact by pulling distance intervals closer by consolidating each weight and minimizing the width of each interval [BMV*22].

Ray entropy: The entropy regularization calculates the entropy of the ray’s distribution [KSH22, BSE*24] for each ray passing through the scene. This involves assessing the uncertainty or randomness in the predicted densities along the ray and avoiding the floaters in the rendered scene. The entropy regularization is described as per Equation 5

$$\mathcal{L}_{ent} = \left(- \sum_{i=1}^N p(r_i) \log(p(r_i)) \right) \quad (5)$$

Here, N is the number of sampled points on the ray r , and p_i is the opacity of each sampled point.

Complete loss. To train the C^3 -NeRF method, we used the combined loss formulation, as described in Equation 6.

$$\mathcal{L}_{total} = \mathcal{L}_{mse} + \lambda_{ent} \mathcal{L}_{ent} + \lambda_{dist} \mathcal{L}_{dist} \quad (6)$$

Here, λ_{ent} and λ_{dist} are the weights for the regularization. We keep these parameters for training the complete network as $\lambda_{ent} = 1e - 3$ and $\lambda_{dist} = 1e - 2$.

4. Experiments

In this section, we perform an exhaustive set of experiments to demonstrate the efficacy of the proposed framework. We start by discussing the datasets, evaluation metrics, and comparison baselines.

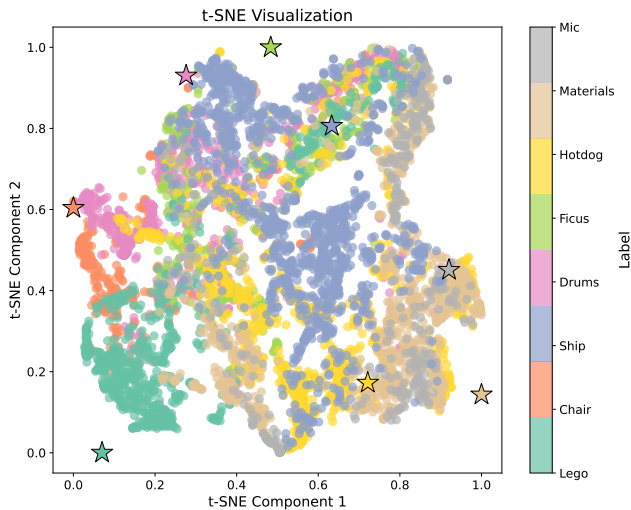


Figure 9: The figure illustrates the representation space of C^3 -NeRF, visualized through a t-SNE map over scenes from NeRF Synthetic 360° dataset [MST*21].

4.1. Datasets and Evaluation Metrics

We evaluate the model performance on widely used NeRF Synthetic 360° [MST*21], Tanks and Temples [KPZK17], and Forward-facing LLFF [MSOC*19] datasets. NeRF Synthetic 360° consists of eight heterogeneous scenes, not necessarily with the same camera parameters. Each scene contains 100 training views and 200 test views of resolution 800×800 rendered from either upper-hemisphere or full-hemisphere. The Real forward-facing LLFF dataset consists of 8 different scenes, each captured using a handheld camera in a forward-facing manner, simulating real-world camera movement. Each scene ranges from 20 to 62 images with 1008×756 resolution. The Tanks and Temples dataset comprises 5 large-scale scenes with complex geometries and real-world objects. Each scene is of size 1920 resolution. We have used the mask version of the tanks and temple dataset for training and testing methods based on the work [LGZL*20]. Besides, we create another dataset with 22 scenes using BlenderNeRF [Raa]. For this, we have used freely available 3D object meshes online, and each scene has a 100 training view and a 100 test view rendered from the upper hemisphere. The purpose of creating this dataset is to stress test the continual learning setup of C^3 -NeRF and empirically analyze the representative upper bound of the network parameters. Moreover, each scene is rendered with the same camera parameters to best adjudge the performance in a continual setup.

Metrics. We quantitatively evaluate the model performance Peak-Signal-to-Noise-Ratio (PSNR) [psn87], Structural Similarity Index (SSIM) [WBSS04], and Perceptual Score (LPIPS) [ZIE*18].

4.2. Training Details

For training C^3 -NeRF, we used a batch size of 10,000 rays and an initial learning rate of $2e-3$. C^3 -NeRF was learned for 30 epochs across both synthetic and real datasets. The model was

Method	Finetuning	NeRF Synthetic 360°		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
NeRF*	x	21.75	0.84	0.16
IBRNet [WWG*21]	\checkmark	28.14	0.94	0.07
MVSNeRF [CXZ*21]	\checkmark	27.07	0.93	0.17
CP-NeRF [HLX*23]	x	29.54	0.92	0.09
C^3 -NeRF (ours)	x	29.40	0.94	0.09
CP-NeRF [HLX*23]	\checkmark	31.77	0.95	0.06
C^3 -NeRF (ours)	\checkmark	32.08	0.95	0.06

Table 1: Quantitative comparison on NeRF Synthetic 360° dataset [MST*21] across different baseline methods. Our conditional-cum-continual framework shows competitive performance with multi-scene modeling and even outperforms the previous state-of-the-art methods. NeRF* indicates vanilla NeRF trained under the proposed conditional-cum-continual setup.

conditioned on scene-specific pseudo labels, enabling it to learn distinctive representations for each scene. The hash table size for Instant-NGP was set to $T = 2^{19}$, with a feature size of $F = 4$. In the case of the NeRF Synthetic 360° dataset, we employed the generative replay method [SLKK17] by rendering previous scenes using the new scene’s camera parameters. For real-world datasets like LLFF [MSOC*19] and Tanks and Temples [KPZK17], which do not provide 360° views, we stored the camera parameters for each scene. This approach allowed us to synthesize novel views without causing artifacts such as floaters or ghosting, which are more easily avoided in the synthetic dataset with full 360° coverage.

4.3. Choice of Comparison Baselines

Since there are no existing recognized works that model multiple scenes in NeRF, we strategically choose and compare our performance with the most closely related baselines that, in a certain limited sense, fit into our experimental setup [CXZ*21, WWG*21, HLX*23, WWW*24]. While SCARF [WWW*24] is closest and concurrent to our work, it accommodates multiple scenes by increasing the number of parameters. It uses a global parameter generator to extract per-scene features. Moreover, we cannot compare the performance with SCARF due to the unavailability of their code.

MVSNeRF [CXZ*21], IBRNet [WWG*21], and CP-NeRF [HLX*23] aim to achieve generalization of NeRF across unseen scenes by using different conditioning mechanisms either through CNN based feature extractor [WWG*21] or 3D cost volume feature [CXZ*21] or meta-learning another hyper-network to generate weights of NeRF MLP [HLX*23]. Since scene generalization can also be viewed as the network’s ability to model multiple unseen scenes via neural featured-based scene conditioning, we use them as our comparison baselines. As we shall see, these methods show improved performance only when fine-tuned over images of the desired test scene but continue to suffer from drastic forgetting. Furthermore, we adapt the proposed conditional-cum-continual paradigm to vanilla NeRF [MST*21] (referred to NeRF*) for a fair comparison. While we compare MVSNeRF [CXZ*21], IBRNet [WWG*21], and NeRF*, both qualitatively and quantita-

Method	Finetuning	Real Forward-Facing		
		PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
IBRNet [WWG*21]	✓	26.73	0.85	0.18
MVSNeRF [CXZ*21]	✓	25.45	0.88	0.19
CP-NeRF [HLX*23]	x	25.41	0.77	0.20
NeRF [MST*21]	per-scene opt.	26.50	0.81	0.25
Instant-NGP [MESK22]	per scene opt.	24.98	0.78	0.24
C^3 -NeRF (ours)	x	22.12	0.66	0.37
CP-NeRF [HLX*23]	✓	27.23	0.81	0.14
C^3 -NeRF (ours)	✓	25.19	0.76	0.25

Table 2: Quantitative comparison on Real Forward Facing dataset [MSOC*19] across different baseline methods. We compare C^3 -NeRF with Instant-NGP (highlighted in red) and vanilla NeRF. The performance of C^3 -NeRF is constrained by that of the Instant-NGP backbone, whose performance falls short of vanilla NeRF itself on the LLFF dataset. However, additional fine-tuning seems to boost C^3 -NeRF’s performance.

Method	Tanks and Temples	
	PSNR \uparrow	SSIM \uparrow
EWC [KPR*17] + NeRF	15.64	0.420
PackNet [ML18] + NeRF	16.71	0.547
MEIL-NeRF [CLBL22]	17.98	0.580
CLNeRF [SLKK17]	21.30	0.640
SCARF [WW*24]	26.78	0.89
C^3 -NeRF (ours)	26.93	0.87

Table 3: Quantitative comparison of C^3 -NeRF on the Tanks and Temple dataset [KPZK17] with baselines that combine either conventional continual learning methods with NeRF or introduce continual learning for single scene optimization along with a concurrent work.

tively, we could only compare CP-NeRF [HLX*23] quantitatively due to the unavailability of their trained weights.

4.4. Experimental Setup

The baseline methods have reported their performance under two experimental settings, i.e., *with* and *without finetuning*.

(a) *Without finetuning*: Given a model trained on a set of scenes, evaluation is done on test scenes without disturbing the learned model parameters.

(b) *With finetuning*: Given a model trained on a set of scenes, the model weights are again fine-tuned over training images of a specific scene for certain epochs before evaluation on its unseen views.

Therefore, we identify the best setting for each baseline for a fair comparison with C^3 -NeRF. Since MVSNeRF [CXZ*21] and IBRNet [WWG*21] are not trained on NeRF Synthetic [MST*21], Real forward-facing [MSOC*19], scenes from these datasets are completely unknown to them. However, NeRF*, CP-NeRF [HLX*23], and C^3 -NeRF are optimized over the scenes from these datasets. Therefore, for a fair comparison we compare the results of MVSNeRF [CXZ*21] and IBRNet [WWG*21] finetuned over these scenes and that of NeRF*, CP-NeRF [HLX*23], C^3 -NeRF without any additional finetuning.

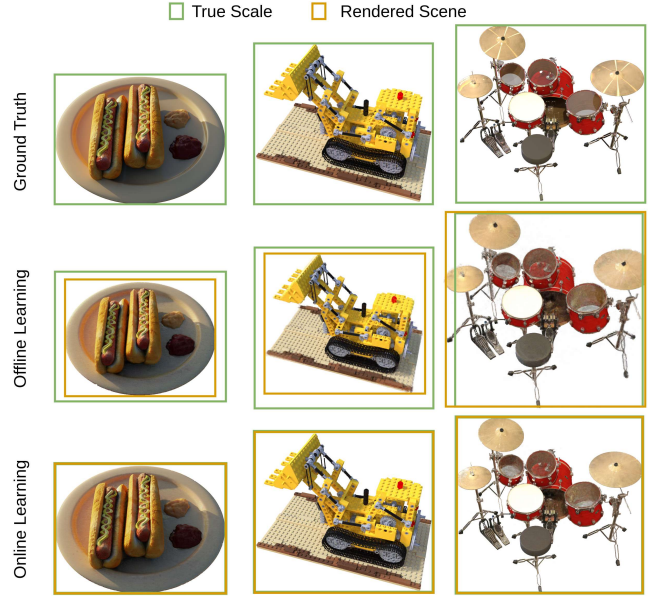


Figure 10: Effect of difference in camera parameters across scenes in the offline setting

Moreover, since CP-NeRF [HLX*23] demonstrates their performance also by finetuning on a specific scene(s), we also report the metrics upon additional finetuning over the respective scene(s). While we understand that this setup is slightly redundant, given the model has already observed those scenes and would obviously enhance the performance, we introduce it here for a fair comparison. Interestingly, if one has access to the training data of certain scenes, which have already been observed and modeled by C^3 during the continual training, finetuning it over the images of such scenes can further boost their rendering quality. As we shall see further, although finetuning can enhance the results, we would like to highlight that it may not always be required since the results without finetuning itself are much more reasonable than the baseline methods.

The summary of the results thus obtained is depicted in Table 1 and Table 2 over the NeRF Synthetic 360° [MST*21] and Real forward-facing [MSOC*19] datasets, respectively.

4.5. Quantitative Analysis

We start by quantitatively comparing C^3 -NeRF with the baselines under the appropriate settings over NeRF Synthetic 360° [MST*21] and Real forward-facing (LLFF) [MSOC*19] datasets, in Table 1 and Table 2, respectively. We observe the best performance over the NeRF Synthetic 360° dataset in terms of SSIM and LPIPS under both settings. However, with no fine-tuning, C^3 -NeRF falls short only by a thin margin (0.14 units of PSNR), especially when conditioning is done without any additional network and the network learns continually. Interestingly, the performance is not the best over the Real forward-facing dataset, as shown in Table 2. The reduced performance is attributed to the limitation of

Method	Current Scene [L]				Current Scene [L -> C]				Current Scene [L -> C -> S]				Current Scene [L -> C -> S -> D]			
	Lego	Chair	Ship	Drums	Lego	Chair	Ship	Drums	Lego	Chair	Ship	Drums	Lego	Chair	Ship	Drums
Offline Sampling	35.62	x	x	x	34.59	34.63	x	x	34.19	34.06	28.53	x	14.41	18.12	28.06	25.02
Online Sampling	35.62	x	x	x	34.81	34.82	x	x	34.24	34.35	28.68	x	33.39	33.89	28.29	25.20

Table 4: Analysing progressive degradation in the rendering quality with offline and online sampling of scenes from the NeRF Synthetic 360° dataset [MST*21]. The cell highlighted in red color shows high degradation in the PSNR value under offline sampling.

Datasets	Offline Sampling			Online Sampling		
	PSNR ↑	SSIM ↑	LPIPS ↓	PSNR ↑	SSIM ↑	LPIPS ↓
NeRF Synthetic 360°	25.211	0.86	0.151	29.40	0.93	0.09
Blender Synthetic	37.735	0.98	0.022	37.734	0.98	0.022

Table 5: Offline vs Online sampling strategy. Analyzing the quantitative performance on training with the available training data (offline) or through images generated through the generative replay (online) while learning C^3 -NeRF. The enhanced performance with online sampling is due to the same camera parameters across all the scenes.

the Instant-NGP-based backbone, which itself is observed to suffer on the forward-facing scenes even when compared to vanilla NeRF as proposed by Mildenhall *et al.* [MST*21] (see Table 2 Row 4 & 5). Furthermore, C^3 -NeRF surpasses Instant-NGP upon finetuning (see Table 2 Row 8). To reiterate, we chose Instant-NGP as the backbone primarily to leverage the advantage of faster training. We even evaluated the standard vanilla NeRF in the proposed conditional-cum-continual setting to demonstrate its inherent ability to accommodate multiple scenes. Given the availability of plenty of NeRF variants, choosing an alternative backbone may eventually perform better, but it is out of the scope of this work at present.

No active baselines directly explore multiple scene modeling in NeRF through the lens of continual learning (except the concurrent work SCARF [WWW*24]). Therefore, as described in SCARF [WWW*24]), we compare C^3 -NeRF with the frameworks that combine either regularization-based Elastic Weight Consolidation (EWC) or parameter isolation - PacketNet [ML18] in continual learning or apply continual learning to a single scene MEIL-NeRF [CLBL22] and CLNeRF [SLKK17] in Table 3 over the Tanks and Temples dataset [KPZK17]. We observe that C^3 -NeRF performs best in terms of PSNR and slightly falls short of SCARF [WWW*24] (just by 0.02 units) in terms of SSIM, even without increasing the number of learnable parameters.

4.6. Qualitative Analysis

We qualitatively demonstrate the visual fidelity of multiple different scenes rendered by C^3 -NeRF across NeRF Synthetic [MSOC*19], real forward-facing [MSOC*19], and Tanks and Temples [KPZK17] datasets in Figure 4, 5, and 6, respectively. Moreover, we perform extensive qualitative comparison across different scenes with MVSNeRF [CXZ*21], IBRNet [WWG*21], and vanilla NeRF [MST*21] in Figure 7 and 8. Overall, we expect

the performance of C^3 -NeRF with multi-scene modeling capability to match that of per-scene optimization-based vanilla NeRF at least closely. We observe that C^3 -NeRF offers the best rendering quality over the NeRF Synthetic dataset and closely matches the vanilla NeRF (per scene optimization) on the Real forward-facing dataset [MSOC*19].

4.7. Online vs. Offline Sampling and Effect of Camera Parameters

While the method does not need to have the training data of the learned scenes by design, it is worthwhile to know how well C^3 -NeRF performs with and without the availability of the actual training data, i.e., offline vs online sampling. Table 5 compares the performance of continually learning a new scene(s) using (a) the training data of the previously encountered scenes - *offline sampling* or using (b) images rendered/sampled from the trained model via the associated scene conditioning - *online sampling*. We find that the performance with *online* and *offline sampling* is nearly the same over the BlenderNeRF dataset, as shown in Table 5, Row 2. Moreover, in Table 5, Row 1, we observe that performance with *online sampling* is, in fact, better than with *offline sampling* over the NeRF Synthetic dataset [MST*21]. We attribute such performance trends to the nature of camera parameters.

Training a neural radiance field under a continual paradigm can be best evaluated when the camera parameters across all scenes are the same because the renderings and associated ground truth images correspond to the same set of camera parameters at every stage of continual training. Interestingly, since the camera parameters of each scene in the BlenderNeRF dataset are the same, it serves as a perfect benchmark for comparing the performance under the two sampling strategies. However, a few scenes in the NeRF Synthetic 360° dataset correspond to different camera parameters. Therefore, the renderings are observed to be misaligned with respect to the ground truth images in the *offline sampling*. Figure 10 shows such misalignments due to differences in camera parameters that eventually cause a reduction in PSNR, leading to reduced performance in the offline sampling. Furthermore, due to differences in camera parameters across scenes, we observe heavier progressive degradation in the rendering quality with *offline* than with the *online sampling*. Table 4 exclusively evaluates the rendering quality of scenes from the NeRF Synthetic dataset with *offline* and *online sampling*.

4.8. Upper bound in C^3 -NeRF

We attempt to explore the upper bound on the number of scenes that can be accommodated within a given set of parameters with an

Method	Per Scene [S, F, H, L]				Fine-tune [S -> F]				Fine-tune [S -> F -> H]				Fine-tune [S -> F -> H -> L]				
	Ship	Ficus	Hotdog	Lego	Ship	Ficus	Hotdog	Lego	Ship	Ficus	Hotdog	Lego	Ship	Ficus	Hotdog	Lego	Avg.
MVSNeRF [CXZ*21]	21.27	19.60	22.44	18.90	12.61	17.72	x	x	15.52	16.83	18.50	x	14.63	16.75	15.80	14.43	15.40
IBRNet [WWG*21]	28.70	29.19	37.14	28.23	22.63	25.23	x	x	23.98	24.63	35.16	x	23.97	24.71	33.49	30.72	28.22
C^3 -NeRF (ours) (no fine-tuning)	30.14	33.95	37.38	35.62	28.86	32.03	x	x	28.56	31.03	35.31	x	28.06	30.61	34.62	32.31	31.4

Table 6: Comparing the extent of information forgetting in MVSNeRF [CXZ*21], IBRNet [WWG*21], and C^3 -NeRF across scenes in the NeRF Synthetic 360° dataset [MST*21].

	Single Dataset		Mixed Dataset	
	PSNR ↑	SSIM ↑	PSNR ↑	SSIM ↑
Groups 1 (8 Scenes)	35.889	0.978	29.40	0.94
Groups 2 (12 Scenes)	33.989	0.971	28.295	0.915
Groups 3 (16 Scenes)	32.448	0.964	27.190	0.906
Groups 4 (22 Scenes)	30.472	0.954	25.744	0.896

Table 7: Analysing upper bound on the number of scenes C^3 -NeRF can accommodate with reasonable loss in rendering quality across single and mixed datasets. Single: 22 scenes from BlenderNeRF dataset. Mixed: 14 scenes from BlenderNeRF + 8 scenes from NeRF Synthetic 360° dataset.

acceptable loss in the rendering quality of the previously learned scenes. Table 7 shows the reduction in the PSNR and SSIM over scenes from a single dataset (Table 7, Column 2) and scenes from a mixed dataset, i.e., scenes from two different datasets (Table 7, Column 3). We used 22 scenes from the BlenderNeRF dataset for the single dataset. Moreover, for the mixed dataset, we combine 8 scenes from the NeRF-Synthetic dataset [MST*21] and 14 scenes from the BlenderNeRF dataset to evaluate cross-dataset performance.

It is pretty evident that increasing the number of scenes for a fixed number of parameters would lead to forgetting. However, the idea here is to understand how accommodating a neural radiance field is. Therefore, we demonstrate that while we do not claim to avoid forgetting completely, most importantly, we succeed in slowing it down to the extent that we can reasonably model around ~ 20 different scenes. We observe that C^3 -NeRF can accommodate over ~ 20 different scenes with PSNR and SSIM higher than or equal to several existing NeRF variants that handle only one scene at a time (compare Table 7 - Row 4, Col 3 and Table 2 - Row 4, Col 3) when evaluated over the scenes from the NeRF Synthetic dataset.

Interestingly, Table 6 compares the extent of forgetting in MVSNeRF [CXZ*21] and IBRNet [WWG*21] with C^3 NeRF (not fine-tune) evaluated over the scenes in the NeRF Synthetic 360° [MST*21] dataset. The experimental setup for the forgetting experiment is that we first fine-tune a scene on MVSNeRF and IBRNet. Later, for each new scene, we use the previous scene’s fine-tuned weights and fine-tune them. For the proposed method, we follow

the continual learning strategy only. No separate fine-tuning is done for C^3 -NeRF.

4.9. Time Analysis

We compare the training, fine-tuning, and rendering times across different NeRF-based methods on the NeRF Synthetic 360° dataset. All experiments were conducted using an NVIDIA RTX Quadro 5000 GPU, ensuring a consistent hardware environment for performance evaluation.

As shown in Table 4.10, NeRF with our proposed conditional-cum-continual method requires approximately 12 days to train on the entire dataset, with rendering times of about 30 seconds per frame, which is relatively slow for practical applications. MVSNeRF [CXZ*21], though not requiring full dataset for training, needs around 1 hour for fine-tuning per scene and has a rendering time of approximately 14 seconds per frame. IBRNet [WWG*21] takes a significantly longer fine-tuning time of around 9 hours per scene and 18 seconds for rendering each frame, highlighting its limitation in time-sensitive scenarios. CP-NeRF [HLX*23], on the other hand, needs 8 RTX 3090 single-day GPU time for training.

In contrast, our proposed C^3 -NeRF method significantly improves time efficiency. It completes training on the full dataset in just around 8 hours, significantly faster than NeRF, its variants, and CP-NeRF. Moreover, fine-tuning time per scene is drastically reduced to approximately 10 minutes, making it suitable for large-scale multi-scene tasks. Most notably, the rendering time per frame is reduced to an impressive 1.2 seconds, greatly outperforming all other methods in terms of rendering efficiency. This makes C^3 -NeRF a practical solution for real-time or near-real-time scene rendering tasks. This substantial reduction in both training and rendering time is one of the key advantages of our proposed work, enabling it to handle complex scenes more efficiently while minimizing computational costs and offering better performance on even higher-end GPU.

4.10. Ablation

We conducted a comprehensive ablation study to identify the most effective method for integrating multi-scene conditioning into the Instant-NGP [MESK22] architecture. Specifically, we explored different configurations for concatenating the pseudo labels and their sinusoidal positional encodings with the scene coordinates and viewing directions. The results of these experiments are summarized in Table 9.

Method	Training Time/ Complete Dataset	Fine-tuning Time/ Scene	Rendering Time/ Frame
NeRF*	~12 days	x	~30 sec
MVSNeRF [CXZ*21]	x	~1 hour	~14 sec
IBRNet [WWG*21]	x	~9 hours	~18 sec
CP-NeRF [HLX*23]	>2 days	x	x
C^3 -NeRF (ours)	~8 hours	~10 mins	~1.2 sec

Table 8: Time Comparison. Training, fine-tuning and rendering time comparison of C^3 -NeRF with other methods on the dataset NeRF Synthetic 360°. All time comparison experiments are done on NVIDIA RTX Quadro 5000 GPU. NeRF* indicates vanilla NeRF trained under the proposed conditional-cum-continual setup.

Experiments	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow
$XYZ \mid C + SE(\theta, \phi)$	24.108	0.812	0.201
$XYZ + \psi(C) + SE(\theta, \phi)$	12.966	0.813	0.460
$XYZ \mid C + \psi(C) + SE(\theta, \phi)$	24.492	0.824	0.184
$XYZ \mid C + \psi(C) + SE(\theta, \phi) + \psi(C)$	25.927	0.880	0.139

Table 9: Ablation on the configuration choice for conditioning Instant-NGP network [MESK22] - the backbone of C^3 -NeRF.

In the baseline configuration, $XYZ \mid C$, we passed the pseudo label alongside the (x, y, z) scene coordinates into the neural hashing module. This configuration represents the most straightforward approach, where the pseudo class label is directly used to differentiate scenes. Furthermore, in $XYZ + \psi(C)$, we only encode scene coordinates with the neural hashing method and concatenate them with the positional encoding of the pseudo label, $\psi(C)$. However, this configuration yielded lower performance across all metrics, suggesting that the pseudo labels are required in neural hashing and scene coordinates to learn sufficient scene discrimination. The configuration $XYZ \mid C + \psi(C)$, where both the raw class label and its sinusoidal encoding are concatenated with the scene coordinates, demonstrated a marked improvement over the baseline. This highlights the importance of combining the class label’s original and encoded forms for better scene representation. In the final and most successful configuration, $XYZ \mid C + \psi(C) + SE(\theta, \phi) + \psi(C)$, we further integrated the spherical encoding of the viewing direction, $SE(\theta, \phi)$, along with an additional instance of the sinusoidal encoding for the class label. This configuration achieved the best results, as shown in Table 9, with the highest PSNR and SSIM scores and the lowest LPIPS score. This setup’s superior performance suggests that combining scene coordinates, viewing direction, and raw and encoded class labels creates a more expressive and robust representation for multi-scene learning. In summary, our ablation study demonstrates that enriching the multi-scene representation with both positional and spherical encodings and the pseudo labels significantly enhances the network’s ability to model multiple scenes and retain the overall performance.

5. Conclusion

In this work, we introduced C^3 -NeRF, a novel conditional-cum-continual framework to encode multiple scenes into the parameters of a single neural radiance field. By leveraging simple pseudo la-

rels for scene conditioning, our approach overcomes the need for complex feature extractors and pre-trained priors, significantly enhancing adaptability to new scenes and reasonably preserving the information of the previously learned scenes. The extensive experimental evaluation demonstrated that C^3 -NeRF excels in both qualitative and quantitative metrics across synthetic and real datasets in learning multi-scene representations. Empirically, we also explored the upper bound on number of scenes that can be modelled by a fixed number of parameters. Although the information forgetting can be minimized further with increase in number of parameters with increasing number of scenes, our goal here was to maximally utilize the representation capacity of a single neural radiance field parameters. We also show the striking comparison of training and inference speed of C^3 -NeRF with the other methods and observe that C^3 -NeRF offers real-time speeds. Interestingly, we show that online sampling under generative replay produces better results compared to offline sampling. While our current generative replay strategy effectively minimizes information forgetting, its efficacy in extremely diverse scene settings remains an area for improvement. Investigating more sophisticated replay techniques or hybrid approaches that combine several continual learning strategies is an interesting future direction. With multi-scene modeling and ability to accommodate ~ 20 scenes, C^3 -NeRF can potentially serve as a valuable 3D asset.

While at present, the proposed framework can be viewed as the one taking first steps towards multi-scene modelling, there are several open research avenues around it. An interesting future work could involve studying if the model learns useful scene priors during a continual learning process to aid in either faster learning over new scenes (i.e., convergence in fewer iterations or higher initial PSNR while learning to model a new scene) or can model new scenes with fewer number of images (few-shot learning). Moreover, dissecting C^3 -NeRF or the Instant-NGP backbone to analyse the kind of scene attributes are handled by the model at each layer to increase physical interpretability is also an interesting avenue to explore. With the current exploration of the representative capacity of NeRFs and the aforementioned future directions, we believe that this work could serve as a primer for a new perspective on designing multi-scene continual neural radiance fields to continue accommodating new scenes without or with minimal loss of information of previously learned scenes.

6. Acknowledgment

This work was supported by the Prime Minister Research Fellowship (PMRF2122-2557) awarded to Prajwal Singh and by the Jibaben Patel Chair in Artificial Intelligence held by Shanmuganathan Raman.

References

- [BDH*23] BAO Y., DING T., HUO J., LI W., LI Y., GAO Y.: Insert-nerf: Instilling generalizability into nerf with hypernet modules. *arXiv preprint arXiv:2308.13897* (2023). 2
- [BLRW16] BROCK A., LIM T., RITCHIE J. M., WESTON N.: Generative and discriminative voxel modeling with convolutional neural networks. *arXiv preprint arXiv:1608.04236* (2016). 4

- [BMV*22] BARRON J. T., MILDENHALL B., VERBIN D., SRINIVASAN P. P., HEDMAN P.: Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5470–5479. [2](#), [8](#)
- [BSE*24] BONOTTO M., SARROCCO L., EVANGELISTA D., IMPEROLI M., PRETTO A.: Combinerf: A combination of regularization techniques for few-shot neural radiance field view synthesis. *arXiv preprint arXiv:2403.14412* (2024). [8](#)
- [BZY*23] BAO C., ZHANG Y., YANG B., FAN T., YANG Z., BAO H., ZHANG G., CUI Z.: Sine: Semantic-driven image-based nerf editing with prior-guided editing field. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 20919–20929. [2](#), [4](#)
- [CCW*23] CHEN Y., CHEN X., WANG X., ZHANG Q., GUO Y., SHAN Y., WANG F.: Local-to-global registration for bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 8264–8273. [5](#)
- [CFHT23] CHEN Z., FUNKHOUSER T., HEDMAN P., TAGLIASACCHI A.: Mobilerf: Exploiting the polygon rasterization pipeline for efficient neural field rendering on mobile architectures. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 16569–16578. [2](#), [4](#)
- [CLBL22] CHUNG J., LEE K., BAIK S., LEE K. M.: Meil-nerf: Memory-efficient incremental learning of neural radiance fields. *arXiv preprint arXiv:2212.08328* (2022). [2](#), [5](#), [10](#), [11](#)
- [CM23] CAI Z., MÜLLER M.: Clnrf: Continual learning meets nerf. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 23185–23194. [2](#), [5](#)
- [CXG*22] CHEN A., XU Z., GEIGER A., YU J., SU H.: Tensorf: Tensorial radiance fields. In *European conference on computer vision* (2022), Springer, pp. 333–350. [2](#), [5](#)
- [CXZ*21] CHEN A., XU Z., ZHAO F., ZHANG X., XIANG F., YU J., SU H.: Mvsnerf: Fast generalizable radiance field reconstruction from multi-view stereo. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 14124–14133. [2](#), [7](#), [8](#), [9](#), [10](#), [11](#), [12](#), [13](#)
- [CYM*23] CHEN J., YI W., MA L., JIA X., LU H.: Gm-nerf: Learning generalizable model-based neural radiance fields from multi-view images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 20648–20658. [2](#), [4](#)
- [DSQ*24] DENG T., SHEN G., QIN T., WANG J., ZHAO W., WANG J., WANG D., CHEN W.: Pglslam: Progressive neural scene representation with local to global bundle adjustment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2024), pp. 19657–19666. [5](#)
- [DTM23] DEBEVEC P. E., TAYLOR C. J., MALIK J.: Modeling and rendering architecture from photographs: A hybrid geometry-and image-based approach. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*. 2023, pp. 465–474. [4](#)
- [FKYT*22] FRIDOVICH-KEIL S., YU A., TANCIK M., CHEN Q., RECHT B., KANAZAWA A.: Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5501–5510. [2](#), [4](#)
- [FXW*23] FANG S., XU W., WANG H., YANG Y., WANG Y., ZHOU S.: One is all: Bridging the gap between neural radiance fields architectures with progressive volume distillation. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2023), vol. 37, pp. 597–605. [2](#), [5](#)
- [GCML23] GORDON C., CHNG S.-F., MACDONALD L., LUCEY S.: On quantizing implicit neural representations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision* (2023), pp. 341–350. [2](#), [5](#)
- [GGSC23] GORTLER S. J., GRZESZCZUK R., SZELISKI R., COHEN M. F.: The lumigraph. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*. 2023, pp. 453–464. [4](#)
- [HLX*23] HE H., LIANG Y., XIAO S., CHEN J., CHEN Y.: Cp-nerf: Conditionally parameterized neural radiance fields for cross-scene novel view synthesis. In *Computer Graphics Forum* (2023), vol. 42, Wiley Online Library, p. e14940. [2](#), [4](#), [9](#), [10](#), [12](#), [13](#)
- [KPR*17] KIRKPATRICK J., PASCANU R., RABINOWITZ N., VENESS J., DESJARDINS G., RUSU A. A., MILAN K., QUAN J., RAMALHO T., GRABSKA-BARWINSKA A., ET AL.: Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences* 114, 13 (2017), 3521–3526. [5](#), [10](#)
- [KPZK17] KNAPITSCH A., PARK J., ZHOU Q.-Y., KOLTUN V.: Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)* 36, 4 (2017), 1–13. [5](#), [9](#), [10](#), [11](#)
- [KSH22] KIM M., SEO S., HAN B.: Infonerf: Ray entropy minimization for few-shot neural volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 12912–12921. [8](#)
- [LCK*24] LEE S., CHOI J., KIM S., KIM I.-J., CHO J.: Few-shot neural radiance fields under unconstrained illumination. In *Proceedings of the AAAI Conference on Artificial Intelligence* (2024), vol. 38, pp. 2938–2946. [2](#), [4](#)
- [LDG18] LIAO Y., DONNE S., GEIGER A.: Deep marching cubes: Learning explicit surface representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (2018), pp. 2916–2925. [4](#)
- [LGZL*20] LIU L., GU J., ZAW LIN K., CHUA T.-S., THEOBALT C.: Neural sparse voxel fields. *Advances in Neural Information Processing Systems* 33 (2020), 15651–15663. [9](#)
- [LH17] LI Z., HOIEM D.: Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence* 40, 12 (2017), 2935–2947. [5](#)
- [LH23] LEVOY M., HANRAHAN P.: Light field rendering. In *Seminal Graphics Papers: Pushing the Boundaries, Volume 2*. 2023, pp. 441–452. [4](#)
- [Max95] MAX N.: Optical models for direct volume rendering. *IEEE Transactions on Visualization and Computer Graphics* 1, 2 (1995), 99–108. [6](#)
- [MESK22] MÜLLER T., EVANS A., SCHIED C., KELLER A.: Instant neural graphics primitives with a multiresolution hash encoding. *ACM transactions on graphics (TOG)* 41, 4 (2022), 1–15. [2](#), [3](#), [4](#), [6](#), [7](#), [10](#), [12](#), [13](#)
- [ML18] MALLYA A., LAZEBNIK S.: Packnet: Adding multiple tasks to a single network by iterative pruning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* (2018), pp. 7765–7773. [5](#), [10](#), [11](#)
- [MON*19] MESCHEDER L., OECHSLE M., NIEMEYER M., NOWOZIN S., GEIGER A.: Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 4460–4470. [4](#)
- [MSOC*19] MILDENHALL B., SRINIVASAN P. P., ORTIZ-CAYON R., KALANTARI N. K., RAMAMOORTHY R., NG R., KAR A.: Local light field fusion: Practical view synthesis with prescriptive sampling guidelines. *ACM Transactions on Graphics (ToG)* 38, 4 (2019), 1–14. [4](#), [8](#), [9](#), [10](#), [11](#)
- [MST*21] MILDENHALL B., SRINIVASAN P. P., TANCIK M., BARRON J. T., RAMAMOORTHY R., NG R.: Nerf: Representing scenes as neural radiance fields for view synthesis. *Communications of the ACM* 65, 1 (2021), 99–106. [1](#), [2](#), [3](#), [4](#), [6](#), [7](#), [8](#), [9](#), [10](#), [11](#), [12](#)
- [NMOG20] NIEMEYER M., MESCHEDER L., OECHSLE M., GEIGER A.: Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2020), pp. 3504–3515. [4](#)
- [PDBW23] PO R., DONG Z., BERGMAN A. W., WETZSTEIN G.: Instant continual learning of neural radiance fields. In *Proceedings of*

- the *IEEE/CVF International Conference on Computer Vision* (2023), pp. 3334–3344. 2, 5
- [PFS*19] PARK J. J., FLORENCE P., STRAUB J., NEWCOMBE R., LOVEGROVE S.: Deepsdf: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2019), pp. 165–174. 2, 4
- [psn87] Implementation of a modified cvsd coder. *International Journal of Electronics* 62, 3 (1987), 473–479. doi:10.1080/00207218708920998. 9
- [Raa] RAAFAT M.: GitHub - maximeraafat/BlenderNeRF: Easy NeRF synthetic dataset creation within Blender — github.com. <https://github.com/maximeraafat/BlenderNeRF>. [Accessed 19-05-2024]. 9
- [RABT17] RANNEN A., ALJUNDI R., BLASCHKO M. B., TUYTELAARS T.: Encoder based lifelong learning. In *Proceedings of the IEEE international conference on computer vision* (2017), pp. 1320–1328. 5
- [RH01] RAMAMOORTHI R., HANRAHAN P.: An efficient representation for irradiance environment maps. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques* (2001), pp. 497–500. 7
- [Rob95] ROBINS A.: Catastrophic forgetting, rehearsal and pseudorehearsal. *Connection Science* 7, 2 (1995), 123–146. 5
- [SG18] STUTZ D., GEIGER A.: Learning 3d shape completion from laser scan data with weak supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 1955–1964. 4
- [SLKK17] SHIN H., LEE J. K., KIM J., KIM J.: Continual learning with deep generative replay. *Advances in neural information processing systems* 30 (2017). 5, 6, 7, 9, 10, 11
- [SLOD21] SUCAR E., LIU S., ORTIZ J., DAVISON A. J.: imap: Implicit mapping and positioning in real-time. In *Proceedings of the IEEE/CVF international conference on computer vision* (2021), pp. 6229–6238. 5
- [SP24] SHIN S., PARK J.: Binary radiance fields. *Advances in neural information processing systems* 36 (2024). 2, 5
- [TCWZ22] TANG J., CHEN X., WANG J., ZENG G.: Compressible-composable nerf via rank-residual decomposition. *Advances in Neural Information Processing Systems* 35 (2022), 14798–14809. 2, 5
- [THM*03] TESCHNER M., HEIDELBERGER B., MÜLLER M., POMERANTES D., GROSS M. H.: Optimized spatial hashing for collision detection of deformable objects. In *Vmv* (2003), vol. 3, pp. 47–54. 6
- [TSM*20] TANCIK M., SRINIVASAN P., MILDENHALL B., FRIDOVICH-KEIL S., RAGHAVAN N., SINGHAL U., RAMAMOORTHI R., BARRON J., NG R.: Fourier features let networks learn high frequency functions in low dimensional domains. *Advances in neural information processing systems* 33 (2020), 7537–7547. 7
- [TTM*22] TEWARI A., THIES J., MILDENHALL B., SRINIVASAN P., TRETSCHK E., YIFAN W., LASSNER C., SITZMANN V., MARTIN-BRUALLA R., LOMBARDI S., ET AL.: Advances in neural rendering. In *Computer Graphics Forum* (2022), vol. 41, Wiley Online Library, pp. 703–735. 4
- [VdMH08] VAN DER MAATEN L., HINTON G.: Visualizing data using t-sne. *JMLR* 9, 11 (2008). 8
- [WBSS04] WANG Z., BOVIK A. C., SHEIKH H. R., SIMONCELLI E. P.: Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612. 9
- [WCLL23] WANG G., CHEN Z., LOY C. C., LIU Z.: Sparsenerf: Distilling depth ranking for few-shot novel view synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (2023), pp. 9065–9076. 2, 4
- [WLL*21] WANG P., LIU L., LIU Y., THEOBALT C., KOMURA T., WANG W.: Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction. *arXiv preprint arXiv:2106.10689* (2021). 2
- [WRB*23] WAN Z., RICHARDT C., BOŽIĆ A., LI C., RENGARAJAN V., NAM S., XIANG X., LI T., ZHU B., RANJAN R., ET AL.: Learning neural duplex radiance fields for real-time view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 8307–8316. 2, 4
- [WT23] WYNN J., TURMUKHAMBETOV D.: Diffusernerf: Regularizing neural radiance fields with denoising diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 4180–4189. 8
- [WWG*21] WANG Q., WANG Z., GENOVA K., SRINIVASAN P. P., ZHOU H., BARRON J. T., MARTIN-BRUALLA R., SNAVELY N., FUNKHOUSER T.: Ibrnet: Learning multi-view image-based rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 4690–4699. 2, 4, 7, 8, 9, 10, 11, 12, 13
- [WWQQ23] WANG Y., WANG J., QU Y., QI Y.: Rip-nerf: learning rotation-invariant point-based neural radiance field for fine-grained editing and compositing. In *Proceedings of the 2023 ACM International Conference on Multimedia Retrieval* (2023), pp. 125–134. 2, 4
- [WWW*24] WANG Y., WANG J., WANG C., DUAN W., BAO Y., QI Y.: Scarf: Scalable continual learning framework for memory-efficient multiple neural radiance fields. *arXiv preprint arXiv:2409.04482* (2024). 5, 9, 10, 11
- [WWX*17] WU J., WANG Y., XUE T., SUN X., FREEMAN B., TENENBAUM J.: Marrnet: 3d shape reconstruction via 2.5 d sketches. *Advances in neural information processing systems* 30 (2017). 4
- [YHL*23] YANG H., HONG L., LI A., HU T., LI Z., LEE G. H., WANG L.: Contranerf: Generalizable neural radiance fields for synthetic-to-real novel view synthesis via contrastive learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2023), pp. 16508–16517. 2, 4
- [YPW23] YANG J., PAVONE M., WANG Y.: Freenerf: Improving few-shot neural rendering with free frequency regularization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (2023), pp. 8254–8263. 2, 4
- [YYTK21] YU A., YE V., TANCIK M., KANAZAWA A.: pixelnerf: Neural radiance fields from one or few images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2021), pp. 4578–4587. 2, 4
- [ZBS*22] ZHANG X., BI S., SUNKAVALLI K., SU H., XU Z.: Nerfusion: Fusing radiance fields for large-scale scene reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2022), pp. 5449–5458. 2, 4
- [ZIE*18] ZHANG R., ISOLA P., EFROS A. A., SHECHTMAN E., WANG O.: The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2018), pp. 586–595. 9
- [ZLCX23] ZHANG L., LI M., CHEN C., XU J.: Il-nerf: Incremental learning for neural radiance fields with camera pose alignment. *arXiv preprint arXiv:2312.05748* (2023). 2, 5