

Divide and Conquer: Rethinking the Training Paradigm of Neural Radiance Fields

Rongkai Ma^{1*†} Leo Lebrat² Rodrigo Santa Cruz² Gil Avraham³
Yan Zuo³ Clinton Fookes⁴ Olivier Salvado²

¹Nvidia, ²CSIRO, Data61, ³Amazon, ⁴Queensland University of Technology

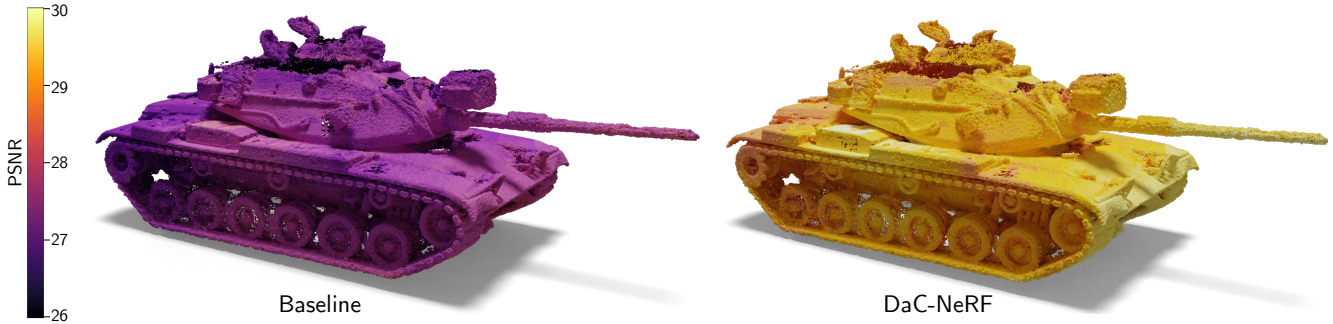


Figure 1. Rendering quality comparison. We conduct a rendering quality evaluation by measuring the Peak Signal to Noise Ratio (PSNR) between rendered views of the specialized models trained using our proposed Divide and Conquer (DaC) training pipeline and one trained through the standard training pipeline which assumes equal importance for all training set images. The PSNR values are projected onto the 3D point cloud of the rendered scene. This comparison highlights that through spatial specialisation, the DaC training paradigm allows for a superior learning of the 3D scene.

Abstract

Neural radiance fields (NeRFs) have exhibited potential in synthesizing high-fidelity views of 3D scenes but the standard training paradigm of NeRF presupposes an equal importance for each image in the training set. This assumption poses a significant challenge for rendering specific views presenting intricate geometries, thereby resulting in suboptimal performance. In this paper, we take a closer look at the implications of the current training paradigm and redesign this for more superior rendering quality by NeRFs. Dividing input views into multiple groups based on their visual similarities and training individual models on each of these groups enables each model to specialize on specific regions without sacrificing speed or efficiency. Subsequently, the knowledge of these specialized models is aggregated into a single entity via a teacher-student distillation paradigm, enabling spatial efficiency for online render-

ing. Empirically, we evaluate our novel training framework on two publicly available datasets, namely NeRF synthetic and Tanks&Temples. Our evaluation demonstrates that our DaC training pipeline enhances the rendering quality of a state-of-the-art baseline model while exhibiting convergence to a superior minimum.

1. Introduction

Rendering high-quality images of real-world scenes remains a significant challenge, particularly those with complex geometries. Neural Radiance Fields (NeRFs [16]) represent a major leap in producing detailed views using volumetric rendering. However, NeRF’s conventional training approach assigns the same importance to all scene perspectives available for training. It uniformly compresses the geometric and photometric information of the scene into the neural network weights. This approach tends to disregard the natural asymmetry of details present in diverse perspectives of complex scenes, leading to declined rendering qual-

*Corresponding author

†Work done while at CSIRO

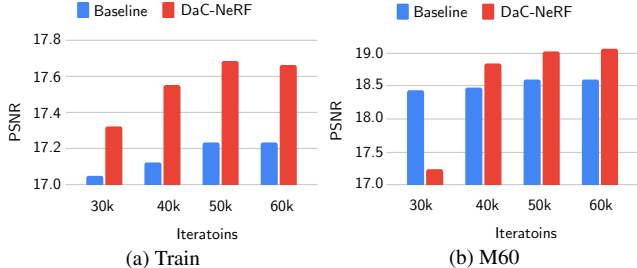


Figure 2. Convergence comparison. We conduct the evaluation of the model trained with our proposed DaC paradigm and the standard pipeline, respectively, on the Train (a) and M60 (b) scenes at 30,000, 40,000, 50,000, and 60,000 iterations. The results demonstrate that the model trained with the DaC’s experts reaches a superior convergence point compared to the standard training strategy.

ity for synthesizing novel views [20]. In a scenario of a flat surface having intricate details on just one side, we expect that perspectives depicting the detailed faces would hold greater significance than those from the uniformly flat side.

Various works have been devised to improve the limitations of NeRFs, including improved space sampling [1–3], explicit spatial feature learning [6, 7, 10], and multi-resolution hashable grid techniques [18] to enhance scene representation. While these methods show improved rendering, they still rely on standard training techniques. We differentiate our work by a novel divide and conquer approach, which is a generic training pipeline and can be applied more universally to any type of NeRF models. We draw inspiration from ensemble learning [5, 23] and mixture of experts (MoE) concepts [22, 28, 30], where multiple models are trained on different scene partitions (known as bootstrapping) and are aggregated during inference. Whilst this technique has seen great success in large-scale NeRF scene reconstruction [25, 30], they fall short when managing computational resources during the inference stage.

Our work. We present a novel training pipeline for neural radiance fields in addressing the complexities of intricate scenes beyond the capabilities of a single NeRF model via a Divide and Conquer (DaC) approach. We initialize our method by dividing the input views into multiple groups based on their inherent structural characteristics. Subsequently, we train an individual expert NeRF model for each group of views, allowing for a more specialized and refined reconstruction of the scene. This can be illustrated in Fig. 1, where the PSNRs of our specialised models (*i.e.*, experts) and the baseline (*i.e.*, NeRF trained on the full dataset) are projected onto the corresponding region of the 3D model. The visualized results distinctly showcase the superior rendering performance of the ensemble of expert models against a single model. Moreover, these expert models are independent and can thus be trained in parallel,

significantly reducing the time complexity of such an approach with parallel training. At the core of our method, we leverage teacher-student distillation [12] to aggregate the knowledge from these experts into one unified model. This approach preserves the original capacity of the NeRF model under consideration without introducing additional time complexity during inference.

Our numerical experiments demonstrate that DaC outperforms the standard NeRF training strategy, particularly in overcoming their novel view rendering performance stagnation. For example, in Figure 2, the performance of the K-Planes NeRF model on the Tanks and Temples dataset tends to plateau after 40k standard training iterations. In contrast, employing the DaC paradigm enables continued improvement in novel view rendering performance beyond this point. Notably, DaC can be applied to most existing NeRF models without introducing runtime or memory overhead during inference. These features make DaC a generic and powerful approach for advancing the next generation of NeRF-based technologies.

2. Related work

The NeRF [16] literature is broad and diverse. However, the works closely related to ours can be categorized into two groups: methods for improving NeRF’s training and rendering efficiency, and approaches involving scene partitioning to create specialized NeRF models for individual partitions. In the following paragraphs, we discuss methods falling within these categories.

Improving speed and quality of NeRFs To improve the training efficiency of a NeRF, a line of research [6, 7, 9–11, 13, 15, 18, 21, 27] propose to replace the coordinate-based multi-layer perceptron adopted in original the NeRF [16] with a sparse (sometimes decomposed) voxel grid of features to represent the radiance field. Specifically, TensorRF [7] propose to model the scene by a factorization of the sparse voxel grid into multiple low-rank components. K-Planes [10] similarly achieves efficiency via decomposing the voxel grid into 3 orthogonal feature planes for static scenes. Thanks to the orthogonal decomposition, K-Planes can represent spatial-temporal features by incorporating an additional temporal axis to model dynamic scenes. However, the trade-off between quality and efficiency is inevitable [13]. Tancik *et al.* [24] propose to mitigate this trade-off by employing a more generalized initialization via meta-learning [8, 16], enabling not only faster convergence but also enhancing rendering quality. While our proposed method can be categorized within the same family, our initialization distinguishes itself by distillation from multiple experts and is not dependent on inner-loop optimization during the inference stage, resulting in faster rendering. Moreover, NeRFLiX [31] enhances rendering quality

by learning the reverse degradation process generated by the radiance field, utilizing two neighbouring views and a NeRF-Style Degradation Simulator (NDS). Another line of research [1–3, 13] is committed to alleviating the aliasing effect induced by NeRF-style sampling via rendering conical frustum instead of rays. Our method can be utilized seamlessly and incorporated into these methods to boost model performance further.

Partitioning Scene of Neural Radiance Field The practice of partitioning large-scale scenes and training a dedicated neural radiance field for each partition is reminiscent of ensemble learning [5]. This approach explored recently [25, 26, 29, 30]. Specifically, BlockNeRF [25] excels in faithfully reconstructing the entire San Francisco city by combining predictions from NeRF models trained on individual city blocks. MegaNeRF [26] leverages a clustering method to partition the scene into a top-down 2D grid which facilitates data parallelism for fast training. SwitchNeRF [30] takes a step further by replacing these heuristic hand-crafted decomposition procedures with an end-to-end learnable gated mixture of experts function, which learns to dispatch a 3D-point to the corresponding specialized NeRF model. Although these methods offer a feasible solution for large-scale scene representation, the substantial memory complexity incurred poses a significant constraint when deploying numerous models into limited computational resources. In contrast, our proposed “divide and conquer” strategy provides a prospective solution for this online memory constraint by consolidating knowledge from multiple specialized models into a unified one using student-teacher distillation [12].

Although our distillation strategy exhibits some resemblances to kiloNeRF [21], the fundamental objective and formulation of our approach diverge markedly. KiloNeRF distills the knowledge of a solitary NeRF trained on the entire scene, directing it towards lower-capacity, space-constrained nerfs to mitigate rendering artifacts in their collaborative rendering process. In contrast, our distillation methodology aims to aggregate knowledge from multiple NeRF experts into a single NeRF model. This consolidated model surpasses the rendering capabilities of a single NeRF trained in the entire scene, producing scene renderings of superior quality.

3. Method

We first provide an overview of our DaC-NeRF to briefly discuss the workflow of our pipeline in Sec. 3.1, which is followed by a detailed discussion of each component of DaC-NeRF (Sec. 3.2).

3.1. Overview

As depicted in Fig. 3, the workflow of our pipeline can be summarized into two stages: divide and conquer. In the first stage, we partition the scene into multiple subsections (refer to the color coding of Fig. 3) such that each scene partition captures specific details of the object. The goal is to ensure that views within each partition share as much relevant information as possible, collectively covering all regions of the object. To achieve this, we design different partitioning methods for object-centric scenes and real-world 360° scenes (see Sec. 3.2.1). Following the partition process, we train a specialized model (*i.e.*, expert) on each of the partitions. This facilitates the local intricacies of the geometry to be learned by each expert. The subsequent distillation stage aggregates the knowledge from each expert into a single unified model; this process not only improves memory efficiency but also facilitates more efficient rendering since a single model will be employed to render new views.

3.2. DaC-NeRF

3.2.1 Dividing Strategy

In this section, we introduce two different scene partition approaches for NeRF datasets. The first approach is designed for object-centric scenes, where cameras are fully described by spherical coordinates. The second approach is tailored for real-world 360-degree scenes, where the positions of the input cameras can be arbitrary.

Dividing object-centric scenes We propose to divide the input camera poses $\mathbf{P} = \{p_1, \dots, p_N\}$ into K even-sized subgroups $\{\mathcal{P}_1, \dots, \mathcal{P}_K\}$ uniquely based on the location of their camera’s centre. While several clustering methods based on the great-circle distance have been explored to divide the set of training views, our empirical findings indicate these methods are highly sensitive to their initial conditions, often leading to uneven clusters. In contrast, we introduce a heuristic that reliably achieves a balanced distribution of cameras across partitions, resulting in a more uniform and coherent division of the scene.

For each camera p_i with associated camera centre $\mathbf{c}_i = (x_i, y_i, z_i)^T \in \mathbb{R}^3$, one can define its *azimuth* coordinate given by $\phi_i = \arctan(x_i, y_i)$. The subsequent step aims at grouping the following N cameras into K distinct subsets $\{\mathcal{P}_1, \dots, \mathcal{P}_K\}$ based on their *azimuthal* angle. Such group allocation is given by,

$$\forall \ell \in \llbracket 1, K \rrbracket, \mathcal{P}_\ell = \left\{ p_k, \frac{2\pi(\ell-1)}{K} \leq \phi_k < \frac{2\pi\ell}{K} \right\}. \quad (1)$$

This heuristic strategy is designed for object-centric scenes where the views’ distribution is approximately uni-

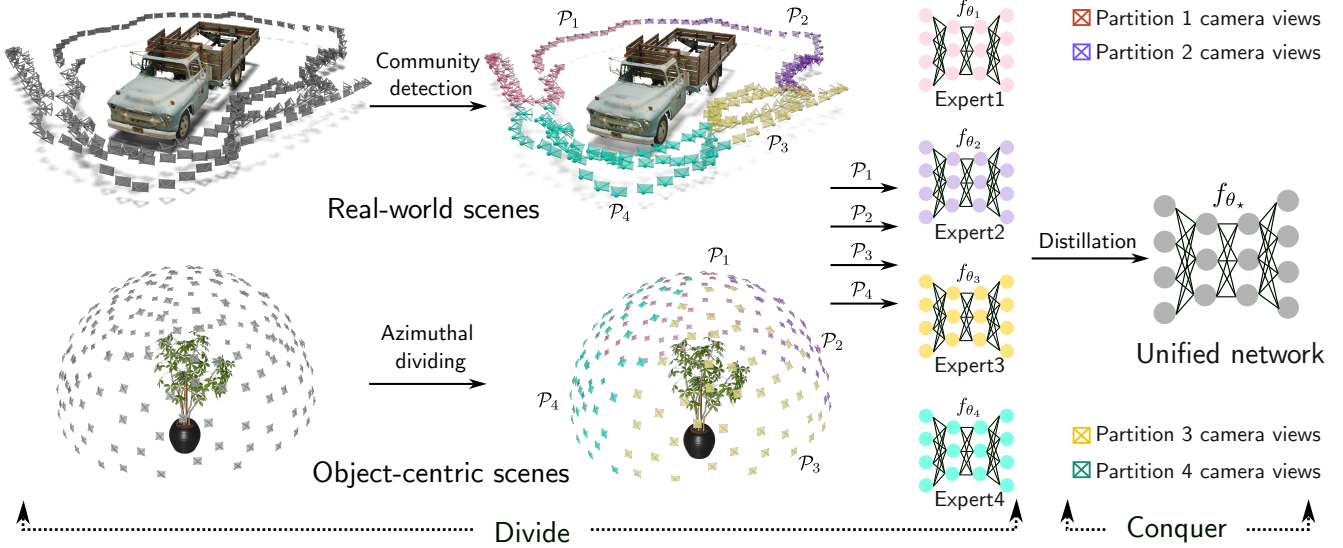


Figure 3. Our Divide and Conquer (DaC) training paradigm. We illustrate our training paradigm with Truck and Ficus. In the dividing stage, we aim to partition the input views $\mathbf{P} = \{p_1, \dots, p_N\}$ into K even-sized groups $\{\mathcal{P}_1, \dots, \mathcal{P}_K\}$. To achieve so, we propose two distinct strategies for dividing object-centric and real-world datasets based. For object-centric scenes, we divide based on the Azimuth angle of input views, while we leverage community detection in Complex Network Analysis [17] to divide for real-world scenes. This is followed by a teacher-student distillation to aggregate the knowledge from expert models into a unified network.

form across the upper half-sphere of the object (refer to Sec. 4.1 for the data generation of the blender scenes).

Dividing real-world scenes To extend this method to free-form capture scenarios, it is necessary to consider extra factors related to both the camera and the scene. These include the orientation of the cameras, the geometry of the scene, and the possibility of occlusion.

It is important to emphasize that the aforementioned strategy is only effective on the fixed object-centric setup where the cameras are chosen to be uniformly distributed on a (half)-sphere. For real-world scenes, where free-form image acquisition is performed, the overlap between two camera frustums can not be described by measuring the distance between their polar angles. To extend this method to free-form capture scenarios, it is necessary to consider extra factors related to both the camera and the scene including the orientation of the cameras, the geometry of the scene, and the possibility of occlusion.

To address this, we reframe this problem through the lens of community detection in Complex Network Analysis [17]. The goal is to group nodes in a graph so that nodes in the same group (community or cluster) are more closely connected to each other than to nodes in other groups. These algorithms often operate on graphs succinctly represented as a weighted adjacency matrix $\mathbf{A} \in \mathbb{N}^{N \times N}$ where A_{ij} denotes the strength of the connection between nodes i and j in the graph. Note that this connection strength is

application-dependent and quantifies the similarity between nodes.

In the context of scene partition, we assign a node to each view and model their connection strength using standard structure from motion (SfM) outputs. Specifically, for any pair of views i and j , we denote their connectivity strength A_{ij} as the total number of 3D points in the SfM’s sparse point cloud that were triangulated from 2D feature correspondence computed on these views. We argue that this measure quantifies the similarity in visual content, field of view, and camera positioning between pairs of views.

Once this adjacency matrix is computed, we can split the input views into K subgroups based on this measure using standard graph community detection algorithms like Louvain community detection [4] or spectral clustering [19]. The former is an efficient ($\mathcal{O}(n \log n)$) heuristic for optimizing graph modularity, while the latter is a clustering technique able to extract clusters with highly non-convex shapes. We encourage the reader to refer to the supplementary material for details of these algorithms.

Training experts Finally, after obtaining K partitions $\mathcal{P}_1, \dots, \mathcal{P}_K$ using either object-centric or real-world scene division strategies, we train K expert NeRF models $f_{\theta_1}, \dots, f_{\theta_K}$ — one for each partition. This process enhances the ability of each expert to learn the local geometric nuances of their assigned scene partition more effectively.

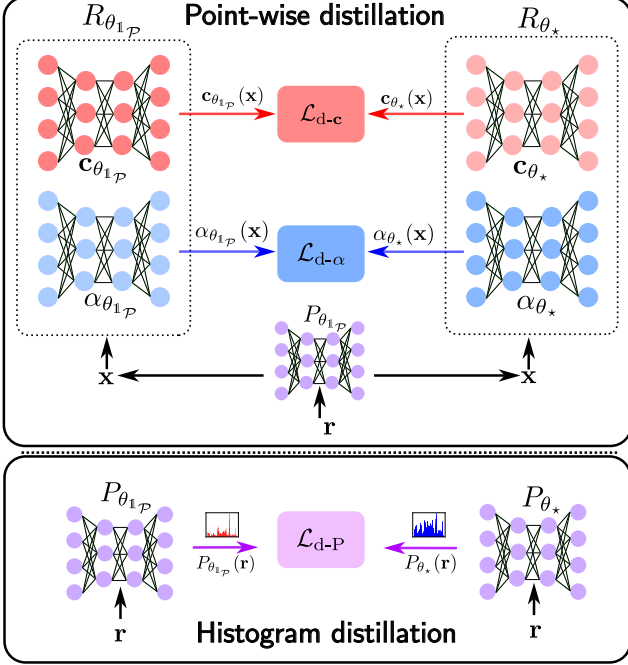


Figure 4. The point-wise distillation loss introduced in our conquer strategy. At the core of our conquer strategy consists of $\mathcal{L}_{d-\alpha}$ to regularize α_{θ_*} , \mathcal{L}_{d-c} to regularize \mathbf{c}_{θ_*} . Additionally, we leverage the histogram loss [2] to guide the learning of P_{θ_*} .

3.2.2 Conquering Strategy

Maintaining K expert models during the inference phase is notably inefficient. To address this, we propose merging the expertise from $f_{\theta_1}, \dots, f_{\theta_K}$ into a single integrated model represented by f_{θ_*} . This integration is achieved through a teacher-student distillation process. Without loss of generality, a NeRF model is composed of a network P_θ that samples points along a ray and a renderer network R_θ to produce the alpha compositing value α and color \mathbf{c} , which are defined by,

$$f_{\theta_i} = (P_{\theta_i}, R_{\theta_i}) \quad (2)$$

$$R_\theta = (\alpha_\theta, \mathbf{c}_\theta), \quad (3)$$

where α_θ and \mathbf{c}_θ are the networks to generate alpha value and color, respectively. We introduce the indicator function to map a 3D point \mathbf{x} along the ray $\mathbf{r}(t)$ into one of the expert subgroups ($f: \mathbb{R}^3 \mapsto \{1, \dots, K\}$) by,

$$\mathbb{1}_{\mathcal{P}_1, \dots, \mathcal{P}_K}(\mathbf{x}) = \left\{ i \mid \exists t, \mathbf{x} \in \mathbf{r}(t) \text{ and } \mathbf{r} \in \mathcal{P}_i \right\}, \quad (4)$$

where $i \in \{1, \dots, K\}$. Likewise, the same process can generalize to mapping rays into one of the expert subgroups by: $\mathbb{1}_{\mathcal{P}_1, \dots, \mathcal{P}_K}(\mathbf{r})$. For simplicity purpose, we denote $\mathbb{1}_{\mathcal{P}_1, \dots, \mathcal{P}_K}(\mathbf{x})$ or $\mathbb{1}_{\mathcal{P}_1, \dots, \mathcal{P}_K}(\mathbf{r})$ as $\mathbb{1}_{\mathcal{P}}$. Having all the

components defined, we can describe the distillation process hereafter.

Given a mini-batch \mathbf{B} of rays $\forall \mathbf{r} \in \mathbf{B}$, $P_{\theta_{1P}}$ first produces a set of refined points $\mathbf{x} \in P_{\theta_{1P}}(\mathbf{r})$ on the rays. Subsequently, we can guide the learning of f_{θ_*} to produce valid opacity α by imposing Mean-Squared Error (MSE) loss on the outputs of α_{θ_*} and $\alpha_{\theta_{1P}}$ (See Fig. 4), which can be described as,

$$\mathcal{L}_{d-\alpha} = \sum_{\mathbf{x} \in P_{\theta_{1P}}(\mathbf{r})} \mathcal{L}_2(\alpha_{\theta_*}(\mathbf{x}), \alpha_{\theta_{1P}}(\mathbf{x})). \quad (5)$$

In a similar fashion, the point-wise MSE loss of the color \mathbf{c} can be obtained as,

$$\mathcal{L}_{d-c} = \sum_{\mathbf{x} \in P_{\theta_{1P}}(\mathbf{r})} \mathcal{L}_2(\mathbf{c}_{\theta_*}(\mathbf{x}), \mathbf{c}_{\theta_{1P}}(\mathbf{x})). \quad (6)$$

Furthermore, to ensure the network P_{θ_*} can produce valid sampled points for R_{θ_*} , we leverage the histogram loss, originally proposed in [2] to align the points generated by P_{θ_*} and $P_{\theta_{1P}}$. This process can be described as,

$$\mathcal{L}_{d-P} = \mathcal{L}_{\text{hist}}(P_{\theta_*}(\mathbf{r}), P_{\theta_{1P}}(\mathbf{r})). \quad (7)$$

Additionally, we keep the loss originally proposed in our baseline implementation and denote them as $\mathcal{L}_{\text{orig}}$ (we encourage readers to refer to the Supplementary Material for the details of $\mathcal{L}_{\text{hist}}$ and $\mathcal{L}_{\text{orig}}$), such that the final loss employed in our conquer phase can be written as,

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{d-\alpha} + \mathcal{L}_{d-c} + \mathcal{L}_{d-P} + \mathcal{L}_{\text{orig}}. \quad (8)$$

It is worth noting that the $\mathcal{L}_{\text{orig}}$ does not include the reconstruction loss between predicted colors and ground truth images since the ground truth images are not required during our conquer phase. This allows f_{θ_*} to assimilate expertise from the experts $f_{\theta_1}, \dots, f_{\theta_K}$, within expanded space beyond the training views (interpolation between training views). Finally, we provide the Pseudo-code of our Divide and Conquer training paradigm in Algorithm 1 and Algorithm 2.

4. Experiments and Discussion

In this section, we provide implementation details of our DaC training paradigm in Sec. 4.1. This is followed by evaluating our method on public datasets in Sec. 4.2. We further conduct ablation studies to justify the design choice of our proposed method in Sec. 4.3.

4.1. Implementation Details

4.1.1 Datasets

We train and evaluate our proposed method using NeRF synthetic datasets [16], including chair, ficus, hotdog, materials and ship. Hereafter, we use NeRF synthetic datasets

Algorithm 1 Divide

Input: A set of posed images $\mathbf{I} = \{I_1, \dots, I_N\}$ **Output:** Expert models $f_{\theta_1}, \dots, f_{\theta_K}$

- 1: **if** Object-centric scenes **then**
 - 2: $\{\mathcal{P}_1, \dots, \mathcal{P}_K\} \leftarrow \text{Azimuth_div}(\mathbf{I})$ §1 Sec. 3.2.1
 - 3: **else if** real-world scenes **then**
 - 4: $\{\mathcal{P}_1, \dots, \mathcal{P}_K\} \leftarrow \text{Graph_cls}(\mathbf{I})$ §2 Sec. 3.2.1
 - 5: **end if**
 - 6: **for** \mathcal{P}_i in $\{\mathcal{P}_1, \dots, \mathcal{P}_K\}$ **do**
 - 7: $f_{\theta_i} \leftarrow \text{Expert_train}(\mathcal{P}_i)$
 - 8: **end for**
-

Algorithm 2 Conquer

Input: Trained experts $f_{\theta_1}, \dots, f_{\theta_K}$ and a batch of rays $\forall \mathbf{r} \in \mathbf{B}$ **Output:** Unified model f_{θ_*}

- 1: $f_{\theta_{i_P}} \leftarrow \text{Exp_retriever}(\mathbf{r})$ Eq. (4)
 - 2: $f_{\theta_*} \leftarrow \text{Distillation}(\mathbf{r}, f_{\theta_{i_P}})$ Eq. (8)
-

and blender datasets interchangeably. Further, we evaluate our method on an unbounded real-world dataset, namely Tanks&Temples [14], which captures 360° of scenes, including M60, Playground, Train, and Truck.

Farthest Point Sampling

We have observed that insufficient quantity of training views in the original blender dataset poses a limitation in training multiple expert models, consequently failing to generalize to novel views. Moreover, the test views of the original blender scenes are distributed non-uniformly on the object, leading to a biased model in evaluation. This motivates us to opt for an alternative way to generate our training and test views for the blender scenes. We utilise Farthest point sampling (FPS) for this task on the blender dataset, which will be discussed briefly hereafter. Given an object centred in $O \in \mathbb{R}^3$, we begin by creating a collection of candidate camera centres $\mathbf{C} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ uniformly distributed across a hemisphere centred at the O with radius r . Subsequently, we select a random point \mathbf{x}_j from \mathbf{C} as the initialization. Thereafter, we compute the distances between \mathbf{x}_j and all other points $\mathbf{x}_{i \neq j} \in \mathbf{C}$, given a predefined distance function f . Finally, we identify the farthest point as the next point of interest and iteratively repeat this process until K points are sampled from our candidate set \mathbf{C} . Due to this being an object-centric dataset, all of the camera directions are pointing towards O , the information of the focal point alone is sufficient to discriminate between the camera’s seen features. Notably, we generate 180 training and 200 test views respectively, across all the scenes in the blender dataset.

4.1.2 Hyperparameters and Loss Functions

We follow the practice in K-planes [10] to train our expert models. Specifically, we train the model for 30k iterations with a learning rate of 0.01. Additionally, we use cosine learning rate scheduler with 512 warm-up iterations. During the conquer stage, we distill the point-wise α value and color \mathbf{c} . However, due to the two-stage ray sampling strategy (proposal-NeRF networks originally proposed in [2]) employed in the baseline model, the distillation process becomes notably intricate. To address this challenge, we begin by acquiring the sampled point locations from the NeRF network of expert models, such that we can perform a forward pass for the NeRF network of f_{θ_*} on these points to get the corresponding α and \mathbf{c} . Then we impose a MSE (Mean Square Error) loss to constrain the disparities between the α and \mathbf{c} parameters of the NeRF network of the experts and those of f_{θ_*} . Another challenge is to regularize the discrepancy of the points sampled by expert models and f_{θ_*} since it is essentially an optimization problem on two histograms with distinct bins. To tackle this challenge, we employ the histogram loss originally imposed between the density weights of NeRF and proposal networks, proposed in [2]. For our specific scenario, we impose the histogram loss between the α values produced by the NeRF network of the expert models and the proposal networks of f_{θ_*} . Note that we use the same learning rate and scheduler as training our expert models. We follow the practice in [21] to further fine-tune the fused model f_{θ_*} with the ground-truth images for 30k iterations.

4.2. Quantitative Results on Public Benchmarks

We compare our proposed DaC against the standard training pipeline on the baseline model (K-Planes [10]) across NeRF synthetic and Tanks&Temples datasets. We report the results on 3 evaluation metrics, including Peak Signal to Noise Ratio (PSNR), Structural Similarity (SSIM), and Multi-Scale Structural Similarity (MS-SSIM).

NeRF Synthetic Based on the findings in Tab. 1, our DaC training paradigm demonstrates improved performance compared to the standard training pipeline across all scenes of the NeRF synthetic dataset. Notably, our DaC consistently outperforms the standard training pipeline on the baseline model, even after an additional 30k iterations of training. This outcome demonstrates the efficacy of our proposed DaC training paradigm for NeRF.

Tanks&Temples In this experiment, we observe a similar finding as noted in the previous section. As Tab. 2 suggests, with our proposed DaC training paradigm the performance of the baseline model is improved across all the 360° unbounded scenes, particularly on the PSNR. This consistent

Scene	Model	PSNR \uparrow	SSIM \uparrow	MS-SSIM \uparrow
Chair	KPlanes	35.05	0.982	0.996
	KPlanes (60k)	35.54	0.983	0.996
	DaC-KPlanes	35.60	0.984	0.996
Ficus	KPlanes	32.02	0.979	0.992
	KPlanes (60k)	32.65	0.982	0.993
	DaC-KPlanes	32.83	0.983	0.993
Materials	KPlanes	31.95	0.966	0.986
	KPlanes (60k)	32.32	0.968	0.990
	DaC-KPlanes	32.47	0.969	0.991
Hotdog	KPlanes	37.67	0.982	0.993
	KPlanes (60k)	38.01	0.983	0.994
	DaC-KPlanes	38.08	0.983	0.994
Ship	KPlanes	30.72	0.888	0.948
	KPlanes (60k)	31.17	0.894	0.952
	DaC-KPlanes	31.29	0.896	0.954

Table 1. The evaluation of our DaC vs. the standard training pipeline on the blender scenes. The results demonstrate consistent improvement over the standard training pipeline across 5 different synthetic scenes (best view in color).

Scene	Model	PSNR \uparrow	SSIM \uparrow	MS-SSIM \uparrow
Truck	KPlanes	21.64	0.663	0.803
	KPlanes (60k)	21.77	0.672	0.811
	DaC-KPlanes	22.01	0.702	0.851
Playground	KPlanes	21.15	0.669	0.823
	KPlanes (60k)	21.30	0.677	0.830
	DaC-KPlanes	21.42	0.678	0.830
Train	KPlanes	17.05	0.574	0.754
	KPlanes (60k)	17.23	0.586	0.766
	DaC-KPlanes	17.66	0.596	0.773
M60	KPlanes	18.43	0.665	0.741
	KPlanes (60k)	18.59	0.872	0.749
	DaC-KPlanes	19.14	0.686	0.772

Table 2. The evaluation of our DaC vs. the standard training pipeline on the Tanks&Temples dataset (best view in color).

enhancement in rendering quality clearly shows the effectiveness of our proposed DaC.

Qualitative Results For qualitative visualisation, we selected rendered test views at 60k iterations from the Tanks&Temples dataset, including M60, Truck, and Train. As illustrated in Fig. 5, the conventional NeRF training pipeline exhibits limitations in the nuances of the captured details within a fixed training budget. In contrast, our proposed DaC clearly improves rendering quality, a distinction particularly evident in the M60 and Truck views, which clearly shows the superiority of our proposed DaC training paradigm.

	Chair	Ficus	Hotdog	Materials	Ship	Avg
2	35.52	32.80	38.01	32.43	31.21	33.39
3	35.61	32.83	38	32.51	31.22	34.03
4	35.60	32.83	38.08	32.47	31.29	34.05
5	35.61	32.82	38.04	32.52	31.24	34.05

Table 3. The analysis of the number of partitions. We conduct the experiments on the NeRF synthetic scenes. Note that the evaluation is reported in PSNR.

4.3. Ablation Study

In this subsection, we conduct various ablation studies to verify our design choices, including investigating the number of partitions, assessing performance with or without overlap between partitions, assessing different partitioning methods for Tanks&Temples, and finally assessing the number of distillation and fine-tuning steps.

4.3.1 The Number of Partitions

In this ablation, we investigate how the number of partitions effect the final result. We conduct experiments on NeRF synthetic scenes with the Peak Signal to Noise Ratio (PSNR) evaluation metric. We compare the test performance across 4 partitioning setups (*i.e.*, 2-5). The results in Tab. 3 suggest that dividing the views into 4 partitions yield the best trade-off between performance and efficiency as increasing the number of partitions over 4 yields no further significant performance gain.

4.3.2 Overlap Between Partitions

We investigate the efficacy of training a dedicated NeRF model on the boundary region between two partitions. This experiment is conducted using the NeRF synthetic scenes, employing the Peak Signal-to-Noise Ratio(PSNR) as the evaluation metric. Specifically, we train a total of six expert models on the boundaries, each covering 120° azimuthal angle with a 60° overlap. The results in Fig. 6 indicate that such an approach does not yield obvious performance improvements, which is consistent across all scenes.

4.3.3 Distillation and Fine-Tuning Steps

Per the findings outlined in Sec. 4.3.1, it is evident that dividing the input views into four groups showcases the best trade-off between local specialization and computational budget. Nevertheless, there remains a necessity to investigate the capacity of how such a training framework could be further enhanced by varying the distillation and fine-tuning with a fixed 60k iterations of training budget. We conduct this experiment on NeRF synthetic scenes using PSNR as the evaluation metric. In Tab. 4, we observe that a balanced

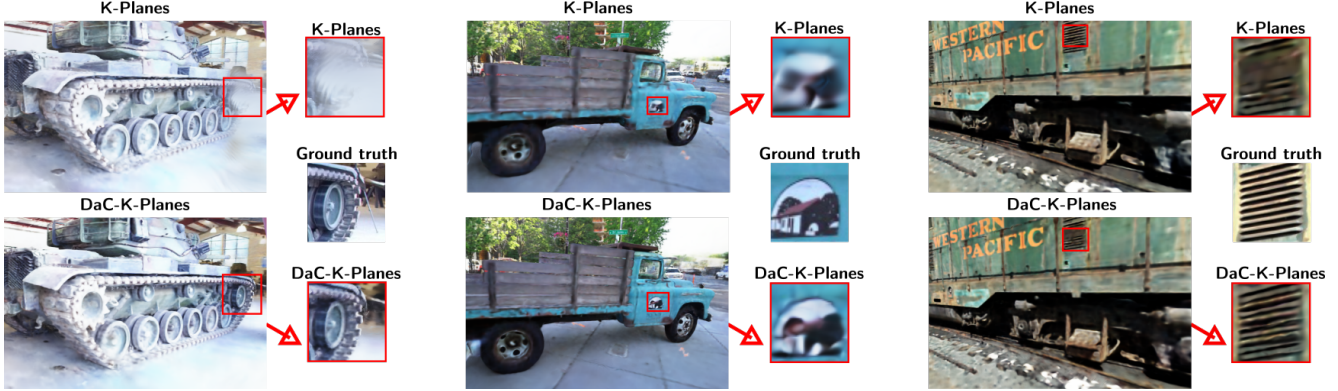


Figure 5. Qualitative Results. At the 60,000 iterations, we select the rendered test views of K-Planes, employing both the conventional training pipeline and our proposed DaC paradigm. The cropped details clearly show that our DaC training paradigm improves the rendering quality for the unbounded real-world scenes.

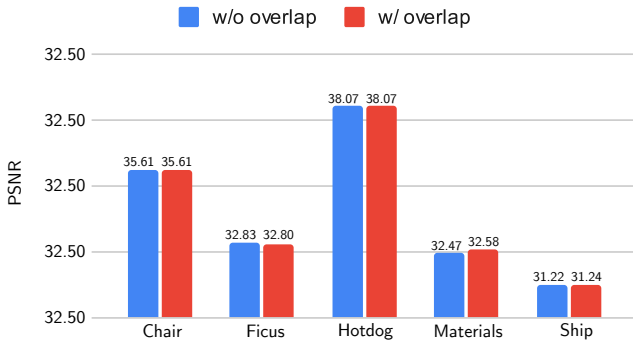


Figure 6. The ablation study on overlapping partition views. In this experiment, we train 6 NeRF models with 60° overlapping azimuth angle. The results demonstrate that training overlapping NeRF models with our DaC yield no significant performance gain.

DST	FT	Chair	Ficus	Hotdog	Materials	Ship	Avg
10k	50k	35.63	32.76	36.10	32.52	30.44	33.49
20k	40k	35.63	32.80	38.00	32.56	31.14	34.03
30k	30k	35.60	32.83	38.08	32.47	31.29	34.05
40k	20k	35.43	32.60	37.87	32.18	31.10	33.83
50k	10k	35.01	31.70	37.16	31.42	30.52	33.16

Table 4. The ablation study of the effect of distillation *vs.* fine-tuning. Note that the evaluation is reported in PSNR. The results clearly suggest that a balanced distillation (DST) and fine-tuning (FT) iterations (30k for each) yields the best result.

iterations between distillation and fine-tuning (denoted as **DST** and **FT** in Tab. 4) yield the optimal performance.

4.3.4 Discussion

Our primary interest focus is to mitigate the limitations inherent in the conventional NeRF training pipeline, specif-

ically in the context of synthetic and unbounded 360° scenes. We demonstrate improved consistent performance over our baseline model for novel view synthesis. The ablation study verifies our design choice.

5. Conclusion and Future Work

In this research, we introduce a pioneering training paradigm for NeRFs, inspired by ensemble based approaches for improving novel view synthesis. Our approach provides insights on how input views of a real-world scene can be effectively partitioned into multiple subgroups, to facilitate the learning of spatial expert models. Our findings demonstrate that these spatial experts achieve superior rendering quality on their local regions. We further showcase the potential of knowledge distillation to aggregate the expertise into a unified model, reducing memory complexity dramatically during the inference stage. Our proposed DaC exhibits frame-work characteristics, making it adaptation-friendly to scene representation methods that may surface in the future.

In future work, we wish to extend such an idea to dynamic scenes. In a similar fashion, one can decompose the time axis into multiple segments and models can learn to specialize on different time partitions such that the final learned model would exhibit better spatial-temporal consistency. Furthermore, another potential extension could be applying our DaC approach into a continual learning setup where the training data is presented in a streaming manner.

References

- [1] 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 *Proceedings of the IEEE/CVF Inter-*

- national Conference on Computer Vision*, pages 5855–5864, 2021. 2, 3
- [2] Jonathan T Barron, Ben Mildenhall, Dor Verbin, Pratul P Srinivasan, and Peter Hedman. Mip-nerf 360: Unbounded anti-aliased neural radiance fields. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5470–5479, 2022. 5, 6
- [3] Jonathan T. Barron, Ben Mildenhall, Dor Verbin, Pratul P. Srinivasan, and Peter Hedman. Zip-nerf: Anti-aliased grid-based neural radiance fields. *ICCV*, 2023. 2, 3
- [4] Vincent D Blondel, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. Fast unfolding of communities in large networks. *Journal of statistical mechanics: theory and experiment*, 2008(10):P10008, 2008. 4
- [5] Leo Breiman. Bagging predictors. *Machine learning*, 24: 123–140, 1996. 2, 3
- [6] Ang Cao and Justin Johnson. Hexplane: A fast representation for dynamic scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 130–141, 2023. 2
- [7] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. Tensorf: Tensorial radiance fields. In *European Conference on Computer Vision*, pages 333–350. Springer, 2022. 2
- [8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, pages 1126–1135. PMLR, 2017. 2
- [9] Sara Fridovich-Keil, Alex Yu, Matthew Tancik, Qinhong Chen, Benjamin Recht, and Angjoo Kanazawa. Plenoxels: Radiance fields without neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5501–5510, 2022. 2
- [10] Sara Fridovich-Keil, Giacomo Meanti, Frederik Rahbæk Warburg, Benjamin Recht, and Angjoo Kanazawa. K-planes: Explicit radiance fields in space, time, and appearance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12479–12488, 2023. 2, 6
- [11] Peter Hedman, Pratul P Srinivasan, Ben Mildenhall, Jonathan T Barron, and Paul Debevec. Baking neural radiance fields for real-time view synthesis. in 2021 ieee. In *CVF International Conference on Computer Vision (ICCV)*, pages 5855–5864, 2021. 2
- [12] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 2, 3
- [13] Wenbo Hu, Yuling Wang, Lin Ma, Bangbang Yang, Lin Gao, Xiao Liu, and Yuewen Ma. Tri-miprf: Tri-mip representation for efficient anti-aliasing neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 19774–19783, 2023. 2, 3
- [14] Arno Knapitsch, Jaesik Park, Qian-Yi Zhou, and Vladlen Koltun. Tanks and temples: Benchmarking large-scale scene reconstruction. *ACM Transactions on Graphics (ToG)*, 36(4):1–13, 2017. 6
- [15] Lingjie Liu, Jiatao Gu, Kyaw Zaw Lin, Tat-Seng Chua, and Christian Theobalt. Neural sparse voxel fields. *Advances in Neural Information Processing Systems*, 33:15651–15663, 2020. 2
- [16] 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. *Computer Vision–ECCV 2020*, pages 405–421, 2020. 1, 2, 5
- [17] EL-MOUSSAOUI Mohamed, Tarik Agouti, Abdessadek Tikniouine, and Mohamed El Adnani. A comprehensive literature review on community detection: Approaches and applications. *Procedia Computer Science*, 151:295–302, 2019. 4
- [18] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multi-resolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 2
- [19] Andrew Ng, Michael Jordan, and Yair Weiss. On spectral clustering: Analysis and an algorithm. *Advances in neural information processing systems*, 14, 2001. 4
- [20] Xuran Pan, Zihang Lai, Shiji Song, and Gao Huang. Activenerf: Learning where to see with uncertainty estimation. In *European Conference on Computer Vision*, pages 230–246. Springer, 2022. 2
- [21] Christian Reiser, Songyou Peng, Yiyi Liao, and Andreas Geiger. Kilonerf: Speeding up neural radiance fields with thousands of tiny mlps. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14335–14345, 2021. 2, 3, 6
- [22] Carlos Riquelme, Joan Puigcerver, Basil Mustafa, Maxim Neumann, Rodolphe Jenatton, André Susano Pinto, Daniel Keysers, and Neil Houlsby. Scaling vision with sparse mixture of experts. *Advances in Neural Information Processing Systems*, 34:8583–8595, 2021. 2
- [23] Douglas R Smith. The design of divide and conquer algorithms. *Science of Computer Programming*, 5:37–58, 1985. 2
- [24] Matthew Tancik, Ben Mildenhall, Terrance Wang, Divi Schmidt, Pratul P Srinivasan, Jonathan T Barron, and Ren Ng. Learned initializations for optimizing coordinate-based neural representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2846–2855, 2021. 2
- [25] Matthew Tancik, Vincent Casser, Xincheng Yan, Sabeek Pradhan, Ben Mildenhall, Pratul P Srinivasan, Jonathan T Barron, and Henrik Kretschmar. Block-nerf: Scalable large scene neural view synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8248–8258, 2022. 2, 3
- [26] Haithem Turki, Deva Ramanan, and Mahadev Satyanarayanan. Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12922–12931, 2022. 3
- [27] Alex Yu, Ruilong Li, Matthew Tancik, Hao Li, Ren Ng, and Angjoo Kanazawa. Plenotrees for real-time rendering of neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5752–5761, 2021. 2

- [28] Seniha Esen Yuksel, Joseph N Wilson, and Paul D Gader. Twenty years of mixture of experts. *IEEE transactions on neural networks and learning systems*, 23(8):1177–1193, 2012. [2](#)
- [29] Yuqi Zhang, Guanying Chen, and Shuguang Cui. Efficient large-scale scene representation with a hybrid of high-resolution grid and plane features. *arXiv preprint arXiv:2303.03003*, 2023. [3](#)
- [30] MI Zhenxing 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*, 2022. [2](#), [3](#)
- [31] Kun Zhou, Wenbo Li, Yi Wang, Tao Hu, Nianjuan Jiang, Xiaoguang Han, and Jiangbo Lu. Nerflix: High-quality neural view synthesis by learning a degradation-driven inter-viewpoint mixer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12363–12374, 2023. [2](#)